

日 本 国 特 許 庁
PATENT OFFICE
JAPANESE GOVERNMENT

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日

Date of Application:

2000年 6月30日

出 願 番 号

Application Number:

特願2000-199738

出 願 人

Applicant(s):

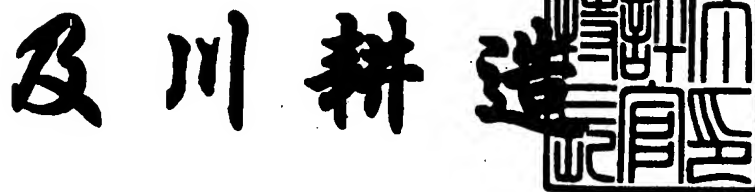
松下電器産業株式会社



CERTIFIED COPY OF
PRIORITY DOCUMENT

2000年12月22日

特許庁長官
Commissioner,
Patent Office



出証番号 出証特2000-3108250

【書類名】 特許願

【整理番号】 2030724007

【提出日】 平成12年 6月30日

【あて先】 特許庁長官殿

【国際特許分類】 G06F 17/27

【発明者】

 【住所又は居所】 大阪府門真市大字門真 1 0 0 6 番地
 松下電器産業株式会社内

 【氏名】 飯塚 泰樹

【特許出願人】

 【識別番号】 000005821

 【氏名又は名称】 松下電器産業株式会社

【代理人】

 【識別番号】 100082692

 【弁理士】

 【氏名又は名称】 蔵合 正博

 【電話番号】 03-5210-2681

【選任した代理人】

 【識別番号】 100081514

 【弁理士】

 【氏名又は名称】 酒井 一

 【電話番号】 03-5210-2681

【手数料の表示】

 【予納台帳番号】 013549

 【納付金額】 21,000円

【先の出願に基づく優先権主張】

 【出願番号】 平成11年特許願第373272号

 【出願日】 平成11年12月28日

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9004843

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 単語分割方式及び装置

【特許請求の範囲】

【請求項 1】 文を単語に分割する単語分割方式において、単語分割されていない文書データから文字結合度を文字間接続確率という形で統計的に計算し、計算した文字間接続確率を分割対象の文に適用し、文字間接続確率が低い部分で分割することで文を単語に分割することを特徴とする単語分割方式。

【請求項 2】 文を単語に分割する単語分割装置において、文の集りである文書データを蓄える文書データ蓄積手段と、単語に分割されていない文書データから文字間接続確率を計算する文字間接続確率計算手段と、計算した文字間接続確率の値を蓄える確率テーブル記憶手段と、計算した文字間接続確率を用いて文を単語単位に分割する文字列分割手段と、文書入力手段と文書出力手段を備えることを特徴とする単語分割装置。

【請求項 3】 文を単語に分割する単語分割方式において、単語分割されていない文書データから文字結合度を文字間接続確率という形で統計的に計算する際、この文字結合度として複数文字からなるある文字列が出現した後にある文字が出現する確率として計算し、この確率を分割対象の文に適用し、文字間接続確率が低い部分で分割することで文を単語に分割することを特徴とする単語分割方式。

【請求項 4】 文を単語に分割する単語分割方式において、単語分割されていない文書データから文字結合度を文字間接続確率という形で統計的に計算する際、この文字結合度として複数文字からなるある文字列が出現した後に複数文字からなるある文字列が出現する確率として計算し、この確率を分割対象の文に適用し、文字間接続確率が低い部分で分割することで文を単語に分割することを特徴とする単語分割方式。

【請求項 5】 文を単語に分割する単語分割方式において、文書データから文字結合度としての文字間接続確率を統計的計算する際、複数文字からなるある文字列が出現した後にある文字が出現する確率と複数文字からなるある文字列が出現する前にある文字が出現する確率とから、複数文字からなるある文字列が出現

した後に複数文字からなるある文字列が出現する確率を計算することを特徴とした前記請求項 4 の単語分割方式。

【請求項 6】 文を単語に分割する単語分割方式において、単語分割されていない学習用文書データから文字結合度としての文字間接続確率を統計的に計算し、計算した文字間接続確率を分割対象の文に適用し、もし分割対象の文の中に文字間接続確率として計算した以外の文字組み合わせが出現した場合は既計算文字間接続確率から目的とする確率を推定し、文字間接続確率が低い部分で分割することで文を単語に分割することを特徴とする単語分割方式。

【請求項 7】 文字間接続確率による単語分割において、分割するかどうかの判断の基準となる確率値の閾値を、分割後の平均単語長をもとに動的に決定することを特徴とする請求項 1 記載の単語分割方式。

【請求項 8】 日本語の文の単語分割において、文字間接続確率の他に、漢字や平仮名といった文字の種類が変化する点で単語が分割されやすいという特徴を併用することで単語分割点を決定することを特徴とした前記請求項 1 の単語分割方式。

【請求項 9】 文字間接続確率の他に、カッコなどの記号部分では必ず単語が分割されるということを併用することで単語分割点を決定することを特徴とした前記請求項 1 の単語分割方式。

【請求項 10】 文を単語に分割する単語分割方式において、単語分割されていない文書データから文字結合度を文字間接続確率という形で統計的に計算して記憶した後、単語辞書を用いて文字列を単語単位に分割する時に単語分割の解の候補が複数あるときは先に計算した文字間接続確率に注目し、分割の解の候補のうち単語分割点における文字間接続確率が小さいものを選択することによって、文を単語に分割することを特徴とする単語分割方式。

【請求項 11】 請求項 10 に記載の単語分割方式において、単語分割の解の候補が複数ある時は、夫々の候補の単語分割点における文字間接続確率の和をその解の候補のスコアとし、分割の解同士をこれらのスコアで比較し、最も低いスコアを取るものを解として選択することを特徴とする単語分割方式。

【請求項 12】 請求項 11 に記載の単語分割方式において、単語分割の解の

候補が複数ある時は、夫々の候補の単語分割点における文字間接続確率の積をその解の候補のスコアとし、分割の解同士をこれらのスコアで比較し、最も低いスコアを取るものを解として選択することを特徴とする単語分割方式。

【請求項 1 3】 請求項 1 0 に記載の単語分割方式において、全ての文字位置についてその文字位置の前が単語分割点だったら文字間接続確率を、単語分割点でなければ定数を与え、それらの確率と定数の和をもって単語分割の解の候補のスコアとし、分割の解の候補同士をこれらのスコアで比較し、最も低いスコアを取るものを解として選択することを特徴とする単語分割方式。

【請求項 1 4】 請求項 1 0 に記載の単語分割方式において、全ての文字位置についてその文字位置の前が単語分割点だったら文字間接続確率を、単語分割点でなければ定数を与え、それらの確率または定数の積をもって単語分割の解の候補のスコアとし、分割の解同士をこれらのスコアで比較し、最も低いスコアを取るものを解として選択することを特徴とする単語分割方式。

【請求項 1 5】 文を単語に分割する単語分割方式において、単語分割されていない文書データから文字結合度を文字間接続確率という形で統計的に計算して記憶した後、単語辞書を用いて文字列を単語単位に分割する時に、辞書に無い文字列を未知語として分割候補に含めることで分割候補を作成し、単語分割の解の候補が複数あるときは先に計算した文字間接続確率に注目し、分割の解の候補のうち単語分割点における文字間接続確率が小さいものを選択することによって、文を単語に分割することを特徴とする単語分割方式。

【請求項 1 6】 請求項 1 5 に記載の単語分割方式において、文字列のある文字位置で終了する単語が辞書にあるがその次の文字位置から開始される単語が辞書に無い時、そこから始まる n 文字以上 m 文字以下の文字列を全て未知語として扱うことで未知語の候補を網羅することを特徴とする単語分割方式。

【請求項 1 7】 請求項 1 5 に記載の単語分割方式において、辞書にある単語には定数スコアを与え、未知語推定された単語にはそれより高い定数のスコアを与え、単語分割の解の候補ごとに単語のスコアの和と分割点における文字間接続確率の和を取ってこれを解の候補のスコアとし、分割の解同士をこれらのスコアで比較し、最も低いスコアを取るものを解として選択することを特徴とする単語

分割方式。

【請求項 1 8】 請求項 1 5 に記載の単語分割方式において、辞書にある単語には定数スコアを与え、未知語推定された単語にはそれより高い定数のスコアを与え、単語分割の解の候補ごとに単語のスコア積と分割点における文字間接続確率の積を取ってこれを解の候補のスコアとし、分割の解同士をこれらのスコアで比較し、最も低いスコアを取るものを解として選択することを特徴とする単語分割方式。

【発明の詳細な説明】

【0 0 0 1】

【発明の属する技術分野】

本発明は、電子計算機を利用した機械翻訳や大量文書検索、テキスト自動要約等を実施する自然言語処理システムの前処理・解析部における方式と装置に関し、特に、文を効率的に単語単位に分割できるようにしたものである。

【0 0 0 2】

【従来の技術】

以後の本発明の説明において、単語は文字の列（文字列）から構成されたもので、文字が組み合わさって意味を形成する単位とする。文（あるいは文章）は単語の列から構成されているものであり、結果として文字列で表される。文書とは文が複数集まってまとまりを作った単位であるとする。

【0 0 0 3】

日本語や中国語など単語を分けて書かない言語を膠着語という。膠着語では、言語の知識を持たない者がその字面だけ見ると、文は長い文字列であって、単語の境界をみつけることができない。

【0 0 0 4】

機械翻訳や自動要約といった自然言語処理システムにおいては、その最初の段階として文の解析が必要になる。日本語のような膠着語では、単語への分割が最初の解析に相当する。

【0 0 0 5】

また文書検索システムでは、例えば「今月の東京都議会」という文字列の「東

京都議会」という語を単語の概念を使わずに文字列検索できるようにしてしまうと、「東京」で検索した場合でも「京都」で検索した場合でもヒットしてしまうことになる。検索語が「京都」だった時に「東京都」がヒットしてしまうのは、本来望まない結果であることから検索ノイズと呼ばれる。こうした検索ノイズを減らすためには、検索対象文書の文を単語に分割しておく必要がある。

【0006】

このような単語分割処理には、通常は辞書を用いた形態素解析処理が使われる。形態素解析では、解析用の辞書を用いて文を単語へ分割するが、形態素解析の精度はこの辞書がどれだけ整っているかに依存する。

【0007】

一方、近年、文書中の文字列や文字の出現といったものを統計的に調べて、処理に必要な情報を得るという提案がなされている。これは例えば、既に単語に分割されている文書から、ある単語（または単語列）の次にどのような単語が出現しやすいかというものを確率として計算し、形態素解析の時にこの情報を使って解の候補を絞るというものである。（参考文献：「単語と辞書」松本祐治 他著、岩波書店 1997年）確率の計算には、単語Nグラムという単語N個組が使われる。NグラムはN-1個の単語の次に、ある単語が出現する確率を計算するもので、この確率計算はマルコフモデルとも呼ばれており、音声認識の単語推定などにも応用されている。ただし単語Nグラムは単語の接続可能性を計算するものであって、辞書にない単語を類推するものではない。また、この単語Nグラムの学習には既に単語に分割されている大量の文書が必要であるが、このような文書は機械的に作ることはできず、人により検査しながら作成する必要があるため、用意するには大きなコストがかかる。

【0008】

この統計処理の考え方をを用いるものとして、単語のNグラムではなく、文字に着目した文字Nグラムがある。文字Nグラムは、(N-1)個の文字列の後にどのような文字が続くかの確率を計算したものである。

【0009】

この文字Nグラムを応用し、文書中に出現する単語になりえそうな文字列の出

現頻度を網羅的に調べて、その文字列前後の文字接続がどれほど散らばっているかを分散という尺度で計算することで単語や慣用句を収集する方法が特開平9-138801に開示されている。

【0010】

論文「正規化頻度による形態素境界の推定」（情報処理学会 自然言語処理研究会 NL-113-3 1996）では、Nグラムの出現頻度を正規化計算することで、辞書を用いずに文を単語単位へ分割する手法が提案されている。

【0011】

辞書を用いない形態素解析は特開平10-326275、特開平10-254874などにも開示されている。特開平10-326275は、文字Nグラムを使って、文字列の部分連鎖確率と単語分割点との関係をテーブルに記憶しておき、そのテーブルを使って単語分割を行うものである。テーブルの作成には、あらかじめ単語単位に分割された文（または文書）を用意し、その文書から計算機により自動学習を行わせる。特開平10-254874も同様に単語単位に分割された文書からあらかじめ学習をする必要がある。

【0012】

【発明が解決しようとする課題】

しかしながら、言語には常に新しい単語が生まれるものであるため、形態素解析用辞書は常にメンテナンスが必要である。また、対象とする文書によって単語の使われ方が違うこともあり、対象とする文書を変更する度に辞書を調整しなければいけない。そして、どれだけ注意していても形態素解析において未知語、すなわち辞書に載っていない単語に遭遇する可能性は否定できず、未知語の出現により形態素解析の精度が低下することがある。

【0013】

辞書を使わないかわりに統計的処理を用いるものとして、前記の通り特開平10-326275などがある。しかしこれらは事前に、単語単位に分割された文書を読ませることでシステムを訓練（自動学習）しておく必要がある。単語単位に分割された文書を用意するためには、人手で文を分割しておくか、または既存の形態素解析システムの出力結果を用いる。だが人手で文を分割するのは多大な

コストが必要であり、文書分野や時代ごとに大量の文を分割して用意することは難しい。そして時代の変化とともに変わっていく言語について、常に大量の分割済み文書を作成し続けなければならない、辞書の整備以上に大変な作業となる。また既存の形態素解析の出力結果を用いた場合、既存の形態素解析における解析失敗部分をそのまま学習してしまい、既存の形態素解析を越える精度は期待できなくなる。

【0014】

本発明は上記の従来技術の課題を解決するためのものであり、基本的に辞書や単語単位に分割された大量の訓練用の文を必要とせず、文を単語へ分割することができる単語分割方式を提供し、また、その方式を実施する装置を提供することを目的としている。

【0015】

【課題を解決するための手段】

この目的を達成するために本発明は、単語に分割されていない文書から文字間の結合度を文字間接続確率という形で計算し、この文字間接続確率を使うことで文を単語単位へ分割するものである。これにより、辞書を使わず、また単語分割後文書を学習する必要もなく、文を単語に分割するという効果を奏するものである。

【0016】

【発明の実施の形態】

（実施の形態1）

以下、本発明の実施の形態について説明する。まず前提となる言語の性質を説明する。文字の出現確率に注目してみる。一般に単語を構成する文字列は、全ての文字の組み合わせの単語が存在するわけではないので、文字の出現は等確率ではない。すなわちある言語の文字の種類をK種類とすると、もし単語を構成する文字が等確率で使われているなら、M文字からなる単語の種類はKのM乗個存在することになる。しかし、実際には語彙数はそれほど多くない。

【0017】

以後の説明では、膠着語として日本語を例にして説明する。日本語の日常生活

で通常使われる文字の種類は、約6千である。この数は、今日の一般的なコンピュータで扱える（JISで規定された）文字の種類数から類推したものである。ここで、日本語2文字の単語について考える。もし全ての文字の組み合わせの単語が存在するなら、6000の2乗＝3千6百万の単語が存在することになる。日本語にはこの他に3文字や4文字の単語も存在するからさらに多くの単語が存在することになる。しかし日本語の総語彙数はたかだか数十万と考えられる。この数は、岩波書店の広辞苑等、日本語辞書の語彙数が20万から30万の間にあることから推測したものである。これについては、「自然言語処理」（長尾真編 岩波書店 1996年出版）の第2章第1節「言語の統計」にも述べられている通りで、一般に文字の出現には偏りがあるものとされている。

【0018】

次に本発明の原理について説明する。ある文字「a」の後に別の文字「b」が続く確率は、もし前記の偏りがなければ、すなわちどんな文字も等確率で出現して単語を形成するなら、言語を構成する文字の種類Kの逆数（日本語の場合約6千分の1）になる。しかし実際には偏りがあるのでそうはならない。具体例で説明する。ある文字列「衆議院」が単語であったとしよう。すると、文字単位での接続確率を「衆」の後に「議」が続く条件付き確率 $P(\cdot\text{議}|\text{衆})$ とした時、この確率は日本語全体を調べてみるなら6千分の1より大きくなるはずである。同様に文字列「衆議」の後に文字「院」が続く条件付き確率 $P(\cdot\text{院}|\text{衆議})$ の場合は前2文字が与えられることから、さらに高い確率を示すはずである。一方で、存在しない単語（文字の組み合わせ）と思われる「衆び」などが出現する確率 $P(\cdot\text{び}|\text{衆})$ は、限りなく0に近くなるはずである。

【0019】

一方、文を構成する単語は、かなり自由な組み合わせが可能である。例えば「これは数学の本だ」「これは音楽の本だ」は両方とも文であるが、「数学」「音楽」の部分は自由な単語が接続できる。つまり単語を構成する文字列「これは」の後に別の単語の文字「数」が続く条件付き確率 $P(\cdot\text{数}|\text{これは})$ は、単語を構成する文字間の接続確率よりも低くなるはずである。この文字間接続確率は、文字の間の結合度と解釈できる。そしてこれが計算できれば、それを基に文字列（

文) を単語単位へ分割することができる。

【 0 0 2 0 】

文字間の接続確率は、文をある程度の量、つまりある程度の量の文書を集めることができれば、そこから統計的に調べて計算することができる。すなわち、文書データベースを構築するような状況ならば、データベースに登録する文書から文字間接続確率を統計的に調べて計算することができる。この計算値は日本語全体について調べた場合の確率値とは違うものであろうが、近似できるものであり、しかもその確率値を調べた文書、あるいは類似の文書の分割に適用するのに適した性質を持つ。

【 0 0 2 1 】

本発明では以上の原理を用いることで、単語分割されていない文書から文字間接続確率を調べ、その文字間接続確率を使うことで、文書を辞書を用いることなく単語単位へと分割する。以下、本発明の実施の形態について図面を用いて説明する。

【 0 0 2 2 】

(実施の形態 1)

図 1 は本発明の実施の形態 1 における単語分割処理方式を説明するフロー図である。図 2 は本実施の形態 1 における文字列分割装置の構成を示すブロック図であり、処理対象文書を電子化された状態で入力するための文書入力手段 2 0 1 と、入力した文書を蓄えておく文書蓄積手段 2 0 2 と、文書から文字間接続確率を計算するための文字間接続確率計算手段 2 0 3 と、計算した確率を記録しておくための確率テーブル格納手段 2 0 4 と、計算した文字間接続確率を使って文書を単語単位に分割するための単語分割手段 2 0 5 と、処理結果の文書を出力する出力手段 2 0 6 を備えている。

【 0 0 2 3 】

以上のように構成された文字列分割装置について、その処理動作を図 1 を用いて説明する。

ステップ 1 0 1 : 文書入力手段 2 0 1 から入力されたデータは、まず文書蓄積手段 2 0 2 に蓄えられる。

ステップ 1 0 2 : このデータから、文字間の接続確率を 2 0 3 の文字間接続確率計算手段が計算し、計算結果を確率テーブル格納手段 2 0 4 に蓄える。計算方法の詳細は後述する。

ステップ 1 0 3 : 文書蓄積手段 2 0 2 に蓄えられたデータについて、確率テーブル格納手段 2 0 4 に蓄えられた確率値を用いることで、文字間の接続確率を調べ、その確率が低い所で分割をし、

ステップ 1 0 4 : 分割された文を出力手段 2 0 4 から出力する。

【 0 0 2 4 】

以上のように本実施の形態 1 における文字列分割装置は、処理対象文書から文字間接続確率を計算し、計算した確率を使って対象文書を単語単位へ分割することができる。

【 0 0 2 5 】

以下、図 1 のステップ 1 0 2 の詳細について説明する。 本発明の第一の実施の形態においては、文字 C_{i-1} と文字 C_i の間の文字間接続確率は、文字列 $C_1 \dots C_{i-1}$ の後に次の文字 C_i が続く条件付き確率で表現することにする。これを式にすると次のように書ける。

【 0 0 2 6 】

【数 1】

$$P(C_i | C_1 C_2 \dots C_{i-1}) \quad (1)$$

【 0 0 2 7 】

しかしこれは計算量が大きく記憶空間が大量に必要なになる。(1) 式のような単語列や文字列の条件付き確率は、一般には N グラム ($N = 1, 2, 3, 4, \dots$)、と呼ばれる文字 N 個組で近似する。文字 N グラムよる条件付き確率とは、その $N - 1$ 個の文字列 $C_{i-N+1} \dots C_{i-1}$ という文字列が続いたという条件のもとで文字 C_i が出現する確率である。すなわち、 N グラムの 1 番目から $N - 1$ 番目の文字列が続いたという条件のもとで N 番目の文字が出現する確率である。こ

れは (2) 式のように書くことができる。

【0028】

【数2】

$$P(C_i | C_{i-N+1} \dots C_{i-1}) \quad (2)$$

【0029】

N グラムの確率は、文字列 $C_1 \cdot C_2 \cdot \dots \cdot C_m$ が調べようとするデータ中に出現する回数を $\text{Count} \cdot (C_1 \cdot C_2 \cdot \dots \cdot C_m)$ とすると、

【0030】

【数3】

$$P(C_i | C_{i-N+1} \dots C_{i-1}) = \frac{\text{Count}(C_{i-N+1} \dots C_i)}{\text{Count}(C_{i-N+1} \dots C_{i-1})} \quad (3)$$

と推定できる（参考文献：「単語と辞書」（松本祐治 他 著 岩波書店 1997年））。

【0031】

なおNグラム計算の時には、計算する文字列（文）の前後にN-1文字の特殊な記号を付与して計算するのが一般的である。（参考文献：同書）これは、一般のNグラムは文の先頭の文字の確率や、文の最後の文字の確率を、特殊記号を含めたNグラムにより計算するからである。N=3の例で説明するなら、「これは本だ」という文字列の3グラムを作成するためには、特殊記号をここでは#で表現することとして、「##これは本だ##」のような文字列を作成してからNグラムを作成する。すると「##こ」「#これ」「これは」「れは本」「は本だ」「本だ#」「だ##」の7つの3グラムを作ることになる。

【0032】

一方、本発明の第1の実施の形態においては一般的なNグラムの計算とは違い

、計算する文字列の前後に $N-1$ 文字の特殊な記号を付与せず、計算する文字列の前にだけ $N-2$ 文字（ただし $N-2 \geq 0$ とし、 $N=1$ の時は 0 とする）の特殊記号を付与する。これは文の後については、文の最後の文字の後には必ず単語として切れるのであるから、文の最後の文字とその後の特殊記号との接続確率を計算する必要がないからである。また文の前については、文の先頭が単語区切であることは自明であるので $N-1$ 個の特殊記号は必要なく、文頭 1 文字と次の 1 文字との接続確率を計算するために特殊記号 $N-2$ 個を含む N グラムを作成する必要がある。

【0033】

先の「これは本だ」の 3 グラムの例で言うなら、文の先頭が単語区切であるのは自明なので、「##こ」によって「##」と「こ」の接続確率を計算する必要は無いが、「#これ」によって「#こ」と「れ」の接続確率を計算する必要がある。文の前に必要な特殊記号の数は $N-2$ 個となる。同様に「本だ#」によって「本だ」と「#」の接続確率を計算する必要はないので、文末に必要な特殊記号の数は 0 個となる。

【0034】

式 (3) を計算し、文字 N 個組とともに計算結果を図 2 の確率テーブル格納手段 204 に記録することが図 1 のステップ 102 に相当する。確率テーブル格納手段 204 は、例えば図 4 (d) のように、 N 文字組とその確率値が格納されるものであるが、文字組で検索しやすく記憶容量も小さくするために、適切な構造を用いて実現されているものとし、ここではその構造を限定しない。

【0035】

ステップ 102 の計算手順は、例えば図 3 に示す手順で実現できる。

ステップ 301：文書を構成する文ごとに、文の前に文頭を表現する特殊記号を $N-2$ 個付与する。

ステップ 302： $N-1$ グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 $N-1$ 個組について、それが何回出現しているかを調べた表を作成する。一般に N グラムの統計を調べる方法は、（参考文献：「言語情報処理」

長尾真 他著、岩波書店 1998 年）などに述べられているが、単純には文

字の種類 K の N 乗を表現できるテーブルを用意し、そこに出現数をカウントして行くか、あるいは文書から全ての N 文字組を取り出しそれをソートして同じものの出現回数をカウントすれば計算できる。

ステップ303： N グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 N 個組について、それが何回出現しているかを調べた表を作成する。これはステップ302と同様である。

ステップ304： N グラム統計の夫々の文字 N 個組文字列について、その出現回数を X とする。同文字 N 個組文字列について、その1番目から $N-1$ 番目の文字列の出現回数をステップ301で作成した $N-1$ グラム統計から調べ、これを Y とする。 X/Y により式(3)の値を計算し、この値を確率テーブル格納手段204に記録する。

【0036】

以上により式(3)の値が計算できるが、 $N-1$ グラムを作成しない方法も存在する。 N グラムは $N-1$ グラム文字列を含むことから、 N グラムを作っておけば $N-1$ グラムの出現頻度も簡単に計算できるからである。

【0037】

以下、文字間接続確率の具体的な計算例を示す。全文書として文字列「a b a a b a」だけが与えられた場合を例とし、ここから文字間接続確率を $N=3$ の N グラム(3グラム)で計算する。まずステップ301として、文(文字列)の前に文頭文末を表現する特殊記号を $N-2$ ($3-2=1$)個付与する。この様子を図4(a)に示す。特殊記号としてここでは#を付けているが、実際には文書に現れない記号を付けるものとする。次にステップ302として、2グラムの統計、すなわち文字2個組の出現回数を調べる。その結果が図4(b)のようになる。同様にステップ303として3グラムの、文字3個組の出現回数を調べ、図4(c)を得る。ステップ304として、図4(b)と図4(c)から、文字3個組についての式(3)の値を計算し、図4(d)を得る。以上が図1のステップ102の詳細説明である。

【0038】

以下では図1のステップ103の詳細を説明する。ステップ103はステップ

1 0 2 で計算した文字間接続確率の表を使って、処理対象の文を構成する文字のそれぞれの部分の接続確率を調べ、分割をする過程である。その計算手順を図 5 に示す。本発明の第一の実施の形態においては、 δ を閾値とし、この値はあらかじめ決められているものとする。

ステップ 5 0 1 : 文書から文を一つ選択する。

ステップ 5 0 2 : 図 3 のステップ 3 0 1 と同様に、文の前に文頭を表現する特殊記号を $N - 1$ 個付与する。

ステップ 5 0 3 : ポインタを文の前に付けた特殊記号の一文字目に移動する。

ステップ 5 0 4 : ポインタ位置から始まる N 文字について、ステップ 1 0 2 で計算した文字間接続確率を調べる。

ステップ 5 0 5 : もしその確率が、あらかじめ決められた閾値 δ 未満だったら、ポインタ位置を 1 文字目とした時の $N - 1$ 文字目と N 文字目の間は単語分割点だったものと推定され、よってそこで分割を行う。もしその確率が閾値 δ 以上だったら、そこは単語分割点ではないので分割を行わない。

ステップ 5 0 6 : ポインタを一文字進める。

ステップ 5 0 7 : ポインタを 1 文字目とした時の N 文字目が文末文字を越える場合、文は終了したものであるとして、ステップ 5 0 8 へ。そうでなければステップ 5 0 4 へジャンプする。

ステップ 5 0 8 : 文書から次の文を選択する。

ステップ 5 0 9 : 次の文が無ければ終了。そうでなければステップ 5 0 2 へ進む。

【 0 0 3 9 】

以上により分割点を発見する。以下、具体的計算例を、先に示した図 4 の文字列「a b a a b a」の $N = 3$ の場合で示す。既に図 4 (d) の文字 3 個組の接続確率は計算されているものとする。閾値として $\delta = 0.7$ が与えられているものとする。まずこの例の場合では文が一つしかないので、ステップ 5 0 1 で「a b a a b a」が選択され、ステップ 5 0 2 で文の前に特殊文字が付けられることで図 4 (a) と同じ状態になる。次にステップ 5 0 3 でポインタを移動させた状態が図 4 (e) である。ここから 3 文字、すなわち「# a b」の確率を図 4 (d)

のテーブルで探すと 1. 0 であり、これは閾値 $\delta = 0. 7$ より大きいので、「# a」と「b」の間は分割されない。以下同様にステップ 5 0 4 からステップ 5 0 6 を繰り返すことで、文字列のそれぞれの点での接続確率が調べられ、この値をもって単語分割点を決定することができる。この様子を図 4 (f) に示す。この例では、単語分割された結果は「a b a / a b a」となる。

【 0 0 4 0 】

もう一つ別の例として、同様に日本語の単純な文字列「にわにわにわ」（庭には二羽）を計算したのが図 6 である。文頭特殊記号を付与したものが図 6 (a)、2 グラムと 3 グラムの出現回数はそれぞれ図 6 (b) (c) に計算され、そこから 3 グラムの文字間接続確率は図 6 (d) のようになる。これを元の文字列にあてはめていくと、図 6 (e) のようになり、閾値 $\delta = 0. 7$ とすることで、結果として「にわ / にわ / にわ」と分割される。

【 0 0 4 1 】

上記 2 例とも文中の文字の種類が少ない例を示した。これは計算例として示すために、非常に短い文から確率計算をしたためである。日本語のように文字種が多い場合は、さらに多くの学習用（確率計算用）の文が必要である。

【 0 0 4 2 】

新聞データ約 1 千万文字からなる文書（文の集合）で文字間の接続確率を計算した例について、その一部を示したものが図 7 である。この計算結果を使い、文字列「利用者の減少と反比例するように」を計算した結果が図 8 である。図 8 では閾値 $\delta = 0. 0 7$ で分割点を決定している。

【 0 0 4 3 】

本発明の第一の実施の形態においては、分割処理対象文書自身から文字間接続確率を計算し、その確率を使って同じ分割処理対象文書を分割した。この方法は対象文書に出現する文字の組み合わせの確率全てを計算できるという点で合理的である。

【 0 0 4 4 】

なお、本発明は、処理対象文書自身だけから文字間接続確率を計算するというものに限定されるものではない。まとまった文書から文字間接続確率を計算して

おき、それを使って別の文書を分割することも可能である。これは漸増的な文書データベースにおいて有効である。この場合は分割対象文書に出現する文字の組み合わせが、確率を計算（学習）した文書に出現していない可能性も否定できないが、これらはNグラム平滑化の問題として（参考文献：「単語と辞書」松本祐治 他著、岩波書店 1997年）などに記述されている方式で対応できる。

【0045】

以上のように、本発明の第1の実施の形態では、ステップ101で入力された文書からステップ102で文字間の接続確率を計算し、この接続確率を使ってステップ103で該文書のそれぞれの文字の接続確率を調べることで単語分割を行い、ステップ104で結果を出力することで、辞書を使わない単語分割を行うことが可能になり、その実用的効果は大きい。

【0046】

（実施の形態2）

本発明の第2の実施の形態の文字列分割装置の構成図は本発明の第1の実施の形態の図2と同じものである。また動作の概要は、本発明第1の実施の形態の図2と同じであるが、計算方式として別のものを用い、よって図1のステップ102、およびステップ103の手順が変更されるので、その詳細を説明する。

【0047】

本発明の第1の実施の形態の説明では、文字列接続確率の計算にはNグラムを用いることで、文字列 $C_{i-N+1} \dots C_{i-1}$ が出現したという条件のもとで文字 C_i が出現する確率を使った。（式（2）参照）例えば文中に出現する「a b c d e f」の「a b c」と「d e f」の間の接続確率を計算するために、「a b c」という文字列が出現した場合に次が「d」である確率を使ったのである。これは既存の技術であるNグラム方式を転用して用いたからである。Nグラムはもともと、単語の接続の確からしさや文字の接続の確からしさを計算し、文全体として正しいかを判断するためのものである。または、いままでに出現した単語列や文字列から次の単語や文字を予想するものである。よって本来は、

【0048】

【数 4】

$$\prod_{i=1}^m P(w_i | w_1 w_2 \dots w_{i-1})$$

という確率式で計算するものを

【0049】

【数 5】

$$\prod_{i=1}^m P(w_i | w_{i-N+1} \dots w_{i-1})$$

と近似した場合の、product・記号 Π の中の項であった。つまり全ての項を掛け合わせる形で使うことが前提だったので、文字列 $C_{i-N+1} \dots C_{i-1}$ が出現したという条件のもとで文字 C_i が出現する確率というように、条件の部分が複数文字（文字列）で、その後に特定の 1 つの文字が出現する確率で扱えたのである。

【0050】

しかし本発明は文字間接続確率を、単語内での文字接続か、単語間の文字接続かを判別するために用いる。よって本発明の第 2 の実施の形態では、文字 C_{i-1} と文字 C_i の接続確率を、ある文字列が出現したという条件でのある 1 文字の出現確率で表現するのではなく、ある文字列が出現したという条件でのある文字列の出現確率を計算するという方式にする。

【0051】

形式的に表現するなら、長さ n 個の文字列 $C_{i-n} \dots C_{i-1}$ が出現したという条件のもとで長さ 1 個の文字列 C_i が出現する確率を計算するよりも、長さ n 個の文字列 $C_{i-n} \dots C_{i-1}$ が出現したという条件のもとで長さ m 個の文字列 $C_i \dots C_{i+m-1}$ が出現する確率を計算する。

【0052】

この確率を式（2）と同様に書くならば、次式（4）のようになる。

【0053】

【数6】

$$P(\underbrace{C_i \dots C_{i+m-1}}_{m \text{ 個}} | \underbrace{C_{i-n} \dots C_{i-1}}_{n \text{ 個}}) \quad (4)$$

【0054】

例えば文中に出現する「a b c d e f」の「a b c」と「d e f」の間の接続確率を計算するために、「a b c」という文字列が出現した場合に次が「d e f」である確率を使うのである。これは $n=3$ 、 $m=3$ の例である。式(4)は $m=1$ の場合が本発明の第1の実施の形態に相当する。

【0055】

また第1の実施の形態が、文の先頭側にある文字列から次の文字列への接続確率を求める、つまり前から後へ進む順方向の確率の計算と考える時、式(4)の $n=1$ 、 $m>1$ という条件は、後から前への接続確率の計算に近似でき、逆方向の確率の計算に相当する。例えば文中に出現する「a b c d e f」の「a b c」と「d e f」の間の接続確率を計算するために、 $n=1$ 、 $m=3$ なら、文字「c」の後に文字列「d e f」が出現する確率ということになるが、文字列「d e f」が出現した場合にその前が「c」である確率に近似できる。これは逆方向の文字間接続確率の計算に相当する。ところが式(4)の計算のためには、 $n+m$ グラムの統計を取る必要がある。しかし、 $n \geq 2$ かつ $m \geq 2$ とすると、4グラム（またはそれ以上）の文字組の統計計算が必要になり、非常に大きな記憶空間を必要とする。

【0056】

そこで本発明の第2の実施の形態では、式(4)の計算を次式(5)で近似する方式を提案する。

【0057】

【数 7】

$$P(C_i|C_{i-n}\dots C_{i-1}) \times P(C_{i-1}|C_i\dots C_{i+m-1}) \quad (5)$$

【0 0 5 8】

式 (5) は、n 個の文字列が出現した後に特定の文字が出現する順方向の確率である第 1 項と、m 個の文字列が出現する前に特定の文字が出現する逆方向の確率である第 2 項の積である。項と文字列の関係を図 1 4 に示す。例えば文中に出現する「a b c d e f」の「a b c」と「d e f」の間の接続確率を計算するために、「a b c」の後に「d」が出現する確率（順方向の第 1 項）と、「d e f」が出現した時にその前が「c」である確率（逆方向の第 2 項）の積をとることを意味する。

【0 0 5 9】

式 (5) の確率値は、第 1 項は n + 1 グラム、第 2 項は m + 1 グラムの統計を取ればよく、式 (6) により計算（推定）できる。

【0 0 6 0】

【数 8】

$$\underbrace{\frac{\text{Count}(C_{i-n}\dots C_i)}{\text{Count}(C_{i-n}\dots C_{i-1})}}_{\text{第 1 項}} \times \underbrace{\frac{\text{Count}(C_{i-1}\dots C_{i+m-1})}{\text{Count}(C_i\dots C_{i+m-1})}}_{\text{第 2 項}} \quad (6)$$

【0 0 6 1】

式 (6) を計算し、文字 n + 1 個組、および文字 m + 1 個組とともに計算結果を、図 2 の確率テーブル格納手段 2 0 4 に記録することが、本発明の第 2 の実施の形態における図 1 のステップ 1 0 2 に相当する。よって、図 2 の確率テーブル格納手段 2 0 4 は、文字 n + 1 個組用と文字 m + 1 個組用の 2 つのテーブルを持つことになる。n ≠ m の場合、この計算手順は、例えば図 9 に示す手順で実現で

きる。

ステップ 9 0 1 : 文書を構成する文ごとに前後に文頭文末を表現する特殊記号を、文頭には $n - 2$ 個、文末には $m - 2$ 個付与する。本発明の第 2 の実施の形態では、後から前への接続確率も計算するため、文の後にも特殊記号を付与する必要がある。

ステップ 9 0 2 : n グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 n 個組について、それが何回出現しているかを調べた表を作成する。

ステップ 9 0 3 : $n + 1$ グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 $n + 1$ 個組について、それが何回出現しているかを調べた表を作成する。

ステップ 9 0 4 : $n + 1$ グラム統計の夫々の文字 $n + 1$ 個組文字列について、その出現回数を X とする。同文字 $n + 1$ 個組文字列について、その 1 番目から n 番目の文字列の出現回数をステップ 9 0 2 で作成した n グラム統計から調べ、これを Y とする。 X / Y により式 (6) の第 1 項の値を計算し、計算結果を確率テーブル格納手段 2 0 4 の文字 $n + 1$ 個組 (第 1 項の確率値) の部分に記録する。

ステップ 9 0 5 : m グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 m 個組について、それが何回出現しているかを調べた表を作成する。

ステップ 9 0 6 : $m + 1$ グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 $m + 1$ 個組について、それが何回出現しているかを調べた表を作成する。

ステップ 9 0 7 : $m + 1$ グラム統計の夫々の文字 $m + 1$ 個組文字列について、その出現回数を X とする。同文字 $m + 1$ 個組文字列について、その 2 番目から $m + 1$ 番目の文字列の出現回数をステップ 9 0 5 で作成した m グラム統計から調べ、これを Y とする。 X / Y により式 (6) の第 2 項の値を計算し、計算結果を確率テーブル格納手段 2 0 4 の文字 m 個組 (第 2 項の確率値) の部分に記録する。

【 0 0 6 2 】

以上が $n \neq m$ の場合である。 $n = m$ の場合、図 2 の確率テーブル格納手段 2 0 4 は、文字 n 個組用だけを用意すればよく、その構造は図 1 2 (d) のように n 文字組とその第 1 項の確率、第 2 項の確率を記録するものとなる。そして $n = m$

の場合、計算手順は例えば図 1 0 に示すように図 9 よりも簡略化できる。

ステップ 1 0 0 1 : 文書を構成する文ごとに前後に文頭文末を表現する特殊記号を $n - 2$ 個ずつ付与する。

ステップ 1 0 0 2 : n グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 n 個組について、それが何回出現しているかを調べた表を作成する。

ステップ 1 0 0 3 : $n + 1$ グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 $n + 1$ 個組について、それが何回出現しているかを調べた表を作成する。

ステップ 1 0 0 4 : $n + 1$ グラム統計の夫々の文字 $n + 1$ 個組文字列について、その出現回数を X とする。同文字 N 個組文字列について、その 1 番目から n 番目の文字列の出現回数をステップ 1 0 0 2 で作成した n グラム統計から調べ、これを Y とする。 X / Y により式 (6) の第 1 項の値を計算し、計算結果を確率テーブル格納手段 2 0 4 の確率値第 1 項の部分に記録する。

ステップ 1 0 0 5 : N グラム統計の夫々の文字 $n + 1$ 個組文字列について、その出現回数を X とする。同文字 N 個組文字列について、その 2 番目から $n + 1$ 番目の文字列の出現回数をステップ 1 0 0 2 で作成した n グラム統計から調べ、これを Y とする。 X / Y により式 (6) の第 2 項の値を計算し、計算結果を確率テーブル格納手段 2 0 4 の確率値第 2 項の部分に記録する。

【 0 0 6 3 】

以上により式 (6) の値が計算できる下地が整ったが、まだ式 (6) の値そのものは求めてなく、実際の値は次の分割過程で計算する。以下では図 1 のステップ 1 0 3 の詳細を説明する。ステップ 1 0 3 はステップ 1 0 2 で計算した文字間接続確率の表を使って、処理対象の文を構成する文字のそれぞれの部分の接続確率を計算し、分割をする過程である。その計算手順を、 $n \neq m$ の場合について図 1 1 に示す。

ステップ 1 1 0 1 : 文書から文を一つ選択する。

ステップ 1 1 0 2 : 図 1 0 のステップ 1 0 0 1 と同様に、文の前後に文頭文末を表現する特殊記号を、文頭には $n - 2$ 個、文末には $m - 2$ 個付与する。

ステップ 1 1 0 3 : ポインタを文の前に付けた特殊記号の一文字目に移動する。

ステップ 1 1 0 4 : ポインタ位置から始まる $n + 1$ 文字について、確率テーブル格納手段 2 0 4 の確率値第 1 項から検索し、それをポインタ位置を 1 文字目とした時の n 文字目と $n + 1$ 文字目の間の確率値第 1 項として記録する。ただし文字と文頭文末特殊記号の間の接続確率は 0 とする。

ステップ 1 1 0 5 : 同様にポインタ位置から始まる $m + 1$ 文字について、確率テーブル格納手段 2 0 4 の確率値第 2 項から検索し、それをポインタ位置を 1 文字目とした時の 1 文字目と 2 文字目の間の確率値第 2 項として記録する。ただし文字と文頭文末特殊記号の間の接続確率は 0 とする。

ステップ 1 1 0 6 : ポインタを一文字進める。

ステップ 1 1 0 7 : ポインタが文末の文字を指していたら、文は終了したものとして、ステップ 1 1 0 8 へ進み、そうでなければステップ 1 1 0 4 へジャンプする。

ステップ 1 1 0 8 : 各文字間について、各文字間に記録された確率値第 1 項と確率値第 2 項の積を取り、式 (6) の値を計算する。それがあらかじめ決められた閾値 δ 未満だったら、そこで分割を行う。もしその確率が閾値 δ 以上だったら、そこは単語分割点ではないので分割を行わない。

ステップ 1 1 0 9 : 文書から次の文を選択する。

ステップ 1 1 1 0 : 次の文が無ければ終了。そうでなければステップ 1 1 0 2 へ進む。以上により分割点を発見する。 $n = m$ の場合も、同様である。

【 0 0 6 4 】

以下、具体的計算例を示す。全文書として文字列「仕事は仕事」だけが与えられた場合を例とし、ここから文字間接続確率を $n = m = 2$ の $(n + 1)$ グラム… (3 グラム) で計算する。まずステップ 1 0 0 1 として、文 (文字列) の前後に文頭文末を表現する特殊記号を $n - 2$ ($3 - 2 = 1$) 個ずつ付与する。この様子を図 1 2 (a) に示す。特殊記号としてここでは # を付けているが、実際には文書に現れない記号を付けるものとする。次にステップ 1 0 0 2 として、2 グラムの統計、すなわち文字 2 個組の出現回数を調べる。その結果が図 1 2 (b) のようになる。同様にステップ 1 0 0 3 として 3 グラムの、文字 3 個組の出現回数を

調べ、図12(c)を得る。またステップ1004として、図12(b)と図12(c)から、文字3個組についての式(6)の第1項の値を計算し、図12(d)の第1項の部分を得る。また、ステップ1005として、図12(b)と図12(c)から、文字3個組についての式(6)の第2項の値を計算し、図12(d)の第2項の部分を得る。

【0065】

注意しなければいけないのは、図12(d)は同じ箇所の文字間接続確率を記録したものではない。図12(d)の表の第2行目「仕事は」について、第1項の部分に記録された確率は「仕事」と「は」の接続確率第1項であり、第2項の部分に記録された確率は「仕」と「事は」の接続確率第2項である。以上により確率値のテーブルができたので、次に図11の処理に進む。

【0066】

ステップ1101として「仕事は仕事」が選択され、ステップ1102で前後に特殊文字が付けられることで図12(a)と同じ状態になる。ステップ1103からステップ1105で図12(e)を得る。第2項は文頭特殊記号と文字との間なので0とする。同様にステップ1103から1106までの繰り返りで図12(f)を得る。ステップ1108により各文字間の接続確率が計算され、図12(f)の接続確率の部分を得る。あらかじめ決められた閾値 δ (ここでは閾値 $\delta = 0.6$ とする)未満の部分で分割を行うと、図12(f)に示す通り、「仕事／は／仕事」という分割結果を得る。

【0067】

なお閾値 δ はあらかじめ決められたものとして扱ってきたが、確率の値を計算した後、単語分割結果が望むべく平均単語長を満すように動的に決めてもよい。すなわち、図13に示すように、閾値 δ が大きければ平均単語長は長くなり、閾値 δ を小さくすれば平均単語長は短くなる。分割結果で調整しながら δ を文書ごとに決めるようにすれば、適切な値が取れるようになる。

【0068】

また、閾値 δ は一率として扱ってきたが、何らかの基準により複数設定してもよい。日本語の場合は本来、平仮名部分の平均単語長は漢字部分の平均単語長よ

りも短い傾向にある。これは平仮名が助詞などの一文字の単語を含むからである。その一方で片仮名部分は外国語の発音を表記したものが多いことから平均単語長は長い。よって閾値 δ を文字種（漢字・平仮名・片仮名）により複数設定してもよい。

【 0 0 6 9 】

また、日本語の場合、単語分割点は文字種（漢字・平仮名・片仮名）の変化点にあることが多い。よって文字種の変化点の閾値 δ を他の部分より下げるなどしてより適切な値に調整してもよい。

【 0 0 7 0 】

なお、本発明の実施の形態では文頭や文末は必ず単語分割点であるとして説明してきたが、この他に句読点の前後、カッコや記号の前後も単語の分割点とみなしてよく、その部分の確率計算を省略することが可能である。あるいはNグラム統計作成において、句読点や記号の前後も文の切れ目として計算してよい。すなわち、本発明第2の実施例で用いた「仕事は仕事」の3グラムは、文頭文末特殊記号を付与することで「#仕事は仕事#」という文字列を作って計算した。これが「仕事は、仕事」だった場合、文は2つであるとみなし、「#仕事は#」と・・・「#仕事#」の2つの文について計算するようにしてもよい。また、式(5)で第1項と第2項の積を取ることで確率を計算したが、第1項と第2項の和、あるいは加重平均などにより確率を計算してもよい。

【 0 0 7 1 】

以上のように、本発明の第2の実施の形態では、文字 n 個の後に文字 m 個が続く確率を近似式(6)で計算することで、より正確に、辞書を使わない単語分割を行うことが可能になり、その実用的効果は大きい。

【 0 0 7 2 】

(実施の形態3)

第3の実施の形態の文字列分割装置は、あらかじめ作成されている単語辞書を持ち、この辞書を利用しながら文字列を単語単位に分割するが、その過程で本発明の第1、または第2の実施の形態で用いた文字間接続確率を利用する。

【 0 0 7 3 】

まず原理を説明する。

【0074】

「小田中学校」という文字列を分割することを考える。単語辞書に「学校」・・「小田」「小田中」「中学校」という4つの単語があったとすると、分割候補は図17(a)のように「小田／中学校」と「小田中／学校」の2つが考えられる。辞書だけの情報ではこれら2つの候補から1つを選択するのは難しい。

【0075】

そこで本発明では「文字間接続確率が小さい部分は単語分割点として尤もらしい」として、複数の単語分割解の候補のうち接続確率が小さい所で分割されているものを選択する。図17(a)の例では文字間接続確率が図17(b)のように計算されているとすると、単語の分割点の部分に注目して文字間接続確率 $P2\cdot$ と $P3\cdot$ を比較し、 $P2\cdot$ の方が小さいことから候補(1)の「小田／中学校」を正解として選択する。

【0076】

この原理は長い文字列に対しても同様に適用できる。「大阪市立山田中学校」という文字列を分割することを考える。単語辞書に「学校」「山田」「市立」・・「大阪」「大阪市」「中学」「中学校」「田中」「立山」といった単語が存在する時、分割解の候補は図18(a)に示すように「大阪／市立／山田／中学校」と「大阪市／立山／田中／学校」の2つの候補が考えられる。この2つのどちらを正解とするかについても、辞書だけの情報から選択するのは難しい。そこで先の例と同様にこの解の候補の選択にも文字間接続確率を用いて、接続確率が小さい所で分割されているものを選択する。

【0077】

長い文字列の場合は、一つの分割解候補の中に複数の分割点があることから、これら分割点の確率の和や積から分割解の候補ごとにスコアを求め、スコアの比較から候補を選択する。図18(a)の例では、文字間の接続確率が図18(b)のように計算されているとすると、各候補のスコアを分割点での文字間接続確率の和として計算することで図18(c)を得る。よってスコアの値の小さい候補(1)の「大阪／市立／山田／中学校」を正解とする。

【 0 0 7 8 】

上記の原理にのっとり、本発明の第 3 の実施の形態について、図面を用いてその処理手順を説明する。図 1 5 は本発明の第 3 の実施の形態を示した構成図の一例である。第 3 の実施の形態の文字列分割装置は、処理対象文書を電子化された状態で入力するための文書入力手段 2 0 1 と、入力した文書を蓄えておく文書蓄積手段 2 0 2 と、文書から文字間接続確率を計算するための文字間接続確率計算手段 2 0 3 と、計算した確率を記録しておくための確率テーブル格納手段 2 0 4 と、あらかじめ作成されている単語辞書を記憶するための単語辞書記憶手段 2 0 7 と、単語辞書の内容を用いて分割解の候補を作成するための分割解候補作成手段 2 0 8 と、解の候補の中から文字間接続確率を用いて解を選択するための解選択手段 2 0 9 と、処理結果の文書を出力する出力手段 2 0 6 を備えている。

【 0 0 7 9 】

処理全体の流れを図 1 6 で示す。

ステップ 1 6 0 1 : 文書入力手段 2 0 1 から入力されたデータは、まず文書蓄積手段 2 0 2 に蓄えられる。

ステップ 1 6 0 2 : このデータから、文字間の接続確率を文字間接続確率計算手段 2 0 3 が計算し、計算結果を確率テーブル格納手段 2 0 4 に蓄える。計算方法の詳細は第 1 の実施の形態、または第 2 の実施の形態の方式に依る。

ステップ 1 6 0 3 : 文書蓄積手段 2 0 2 に蓄えられたデータについて、単語辞書記憶手段 2 0 7 に蓄えられた単語辞書情報から分割解候補作成手段 2 0 8 が分割解の候補を作成し、この分割候補の中から確率テーブル格納手段 2 0 4 に蓄えられた確率値を用いることで解選択手段 2 0 9 が候補を選択し、その結果として文字列が単語単位に分割され（詳細は後述）、

ステップ 1 6 0 4 : 分割された文を文書出力手段 2 0 6 から出力する。

【 0 0 8 0 】

以上のように本発明における文字列分割装置は、処理対象文書から文字間接続確率を計算し、計算した確率と単語辞書の情報を使って対象文書を単語単位へ分割することを特徴とする。以下、図 1 6 のステップ 1 6 0 3 の詳細について、図を用いて説明する。図 1 6 のステップ 1 6 0 3 は、図 1 9 で示す処理を行う。

ステップ1901：分割しようとする文字列の最初の文字から順に、そこから始まる単語が単語辞書記憶手段207の中にあるかを調べ、あればそれらを羅列する。文字列「大阪市立山田中学校」を例とすると、この状態は図20に相当する。

ステップ1902：単語を結ぶことで文の最後まで到達するものを解の候補とする。そして夫々の解の候補のスコアを計算する。スコアは、夫々の解の候補における各単語分割点の文字間接続確率の和とする。文字列「大阪市立山田中学校」を例とすると、解の候補は図18(a)、文字間接続確率が図18(b)の通りの時、夫々の解の候補のスコアは図18(c)となる。

ステップ1903：最も小さいスコアを持つものを解として選択する。図18(c)では候補(1)を選択する。

【0081】

以上により単語分割を得る。文字間接続確率は必ず0以上であり、ステップ1902で文字間接続確率の和を取ったことから、文字列の分割について「中／学校」と「中学校」のように一つの単語になるか複数に分割されるかで複数の候補があった場合は、分割点が少ない方が必ず選択される。ステップ1902では分割解の候補のスコアとして文字間接続確率の和を取り、そのスコアが最小になるものを解とした。これを式で表現すると式(7)のようになる。すなわち、単語分割位置の集合をSとする時、Sに含まれる全ての位置iにおける文字間接続確率 P_i の和が最小となるような集合Sを選択するという意味である。

【0082】

【数9】

$$\arg \min_S \sum_{i \in S} P_i \quad (7)$$

【0083】

本発明はスコアの計算は、確率の和に限定されるものではない。例えばスコア

の計算として文字間接続確率の積を取ってもよい。これを式(7)と同様に記述するなら、次の式(8)のようになる。

【0084】

【数10】

$$\arg \min_S \prod_{i \in S} P_i \quad (8)$$

【0085】

確率の積を計算する場合は、同じ効果を得るものとして対数を取ってもよい。対数を取るにより、積の計算が対数の和へ換算することができる。これは式(9)式(10)のように書くことができる。

【0086】

【数11】

$$\arg \min_S \log \left(\prod_{i \in S} P_i \right) \quad (9)$$

$$= \arg \min_S \sum_{i \in S} \log P_i \quad (10)$$

【0087】

なお、本発明は文字間接続確率を導入して単語分割の解の候補のスコア計算と選択の方法を提案とするものであり、解候補を得るために単語を結んでスコアを計算していく手順については、本実施の形態の手順に特定するものではない。一般にステップ1902のような状況で単語を結んで、図18(a)のような解の候補を作成し、同時にスコア(あるいはコスト)を計算する場合、計算手法としては動的計画法などの探索手法が提案されており、この動的計画法のアルゴリズムを用いてもよい。

【 0 0 8 8 】

以上のように、本発明の第 3 の実施の形態においては、単語分割処理対象文書自身から文字間接続確率を計算しておき、分割解の候補の作成には単語辞書を用いるが、複数の解候補から一つを選択するには文字間接続確率が最小になるものを選択する。よって分割候補選択のための知識の学習に、あらかじめ人手で作成された分割正解文書を大量に用意する必要がないことからコストを抑えることが可能であり、分割対象の文書から確率の形で知識を自動学習することが可能で、文書分野に適合した学習ができる点で合理的であり、その実用的効果は大きい。

【 0 0 8 9 】

(実施の形態 4)

以下、本発明の第 4 の実施の形態について、図面を用いて説明する。図 2 1 は本発明の第 4 の実施の形態を示した構成図の一例である。

【 0 0 9 0 】

第 4 の実施の形態の文字列分割装置は、処理対象文書を電子化された状態で入力するための文書入力手段 2 0 1 と、入力した文書を蓄えておく文書蓄積手段 2 0 2 と、文書から文字間接続確率を計算するための文字間接続確率計算手段 2 0 3 と、計算した確率を記録しておくための確率テーブル格納手段 2 0 4 と、あらかじめ作成されている単語辞書を記憶するための単語辞書記憶手段 2 0 7 と、未知語候補を特定する未知語特定手段 2 1 0 と、単語辞書の内容と未知語特定手段の結果を用いて分割解の候補を作成するための分割解候補作成手段 2 0 8 と、解の候補の中から文字間接続確率を用いて解を選択するための解選択手段 2 0 9 と、処理結果の文書を出力する出力手段 2 0 6 を備えている。

【 0 0 9 1 】

第 4 の実施の形態の文字列分割装置の処理全体の流れを図 2 2 で示す。

ステップ 2 2 0 1 : 文書入力手段 2 0 1 から入力されたデータは、まず文書蓄積手段 2 0 2 に蓄えられる。

ステップ 2 2 0 2 : このデータから、文字間の接続確率を文字間接続確率計算手段 2 0 3 が計算し、計算結果を確率テーブル格納手段 2 0 4 に蓄える。計算方法の詳細は第 1 の実施の形態、または第 2 の実施の形態の方式に依る。

ステップ2203：文書蓄積手段202に蓄えられたデータについて、単語辞書記憶手段207に蓄えられた単語辞書情報と未知語特定手段210が特定した未知語候補を使い、分割解候補作成手段208が分割解の候補を作成し、この分割候補の中から確率テーブル格納手段204に蓄えられた確率値を用いることで解選択手段209が候補を選択し、その結果として文字列が単語単位に分割される。

ステップ2204：分割された文を、出力手段206から出力する。

【0092】

以上のように本発明における文字列分割装置は、処理対象文書から文字間接続確率を計算し、計算した確率と単語辞書の情報、および未知語推定手段が推定した未知語候補を使って対象文書を単語単位へ分割することを特徴とする。以下、図22のステップ2203の詳細について、図と例を用いて説明する。図22のステップ2203は、図23で示す処理を行う。例として文字列「大阪市立山田中学校」を用いる。単語辞書記憶手段207には「学校」「市立」「大阪」「大阪市」「中学」「中学校」「田中」「立山」という単語があるが「山田」は単語として登録されていないものとする。（図24）

ステップ2301：分割しようとする文字列の最初の文字から順に、そこから始まる単語が単語辞書記憶手段207の中にあるかを調べ、あればそれらを羅列する。文字列「大阪市立山田中学校」を例とすると、この状態は図25（a）に相当する。図25（a）の状態では「山田」はまだ単語として認められていない。

ステップ2302：ある文字位置 i で、その直前の文字位置 $i-1$ で終了する単語が存在するのに文字位置 i から開始される単語が存在しない時、文字位置 i から始まる長さ n 文字以上 m 文字以下の全ての文字列を未知語として羅列する。図25（a）の例では5番目の文字「山」の直前で終わる単語「市立」が存在するのに、「山」から始まる単語が存在しない。そこで「山」から始まる長さ $n \sim m$ 文字の文字列を未知語として先の羅列に加える。 n として2、 m として3を取った場合、「山田」と「山田中」の2つが未知語として推定される。この様子が図25（b）に相当する。

ステップ2303：単語を結ぶことで文の最後まで到達するものを解の候補とす

る。図25(b)の中から単語を結ぶことで最後まで到達するものは、図26に示すグラフのようになることから、図27に示す3つが解の候補となる。そして夫々の解の候補のスコアを計算する。スコアは、夫々の解の候補における各単語分割点の文字間接続確率の和とする。図27の分割解候補の夫々のスコアは、文字間接続確率が図18(b)の通りとすると、図27の右端に示す通りである。ステップ2304：最も小さいスコアを持つものを解として選択する。図27では候補(1)が選択される。

【0093】

以上により単語分割を得る。ステップ2303では分割解の候補のスコアとして文字間接続確率の和を取ったが、本発明はスコアの計算はこれに限定されるものではない。スコアの計算方法は式(7)の他に、式(8)式(9)式(10)などが考えられる。上記は単語分割の解の候補について、分割部分の文字間接続確率値だけでスコアを計算したが、本発明では分割しない部分についても特定の値を式に加えることを提案する。

【0094】

解候補のうち、単語分割点については文字間接続確率を、単語分割点以外の文字間には定数を与え、それらから候補のスコアを計算する。より形式的に説明するなら、以下のようなになる。文字の位置についての集合を N 、単語分割する文字位置の集合を S とする時($S \subseteq N$)、各文字位置 i について値 Q_i を次のように定める。すなわち、 S に含まれる位置 i については文字間接続確率を、 S に含まれない位置 i については定数値 T_h を与え(式(12))、解の候補についてこの値の和を取ったものをスコアとし、スコアが最小になるものを解として選択する(式(11)あるいは式(13))。定数値 T_h は本発明の第1の実施の形態における閾値 δ に相当する。

【0095】

【数 12】

$$\arg \min_{S \subseteq N} \sum_{i \in N} Q_i \quad (11)$$

$$Q_i = \begin{cases} P_i, & i \in S \\ Th, & i \notin S \end{cases} \quad (12)$$

【0096】

【数 13】

$$\arg \min_{S \subseteq N} \left(\sum_{i \in S} P_i + \sum_{i \notin S} Th \right) \quad (13)$$

【0097】

例として文字列「新宿泊棟」（新しい宿泊用建物の意味）の分割過程を見る。
「新宿泊棟」の分割候補として図 28（a）のように「新／宿泊／棟」と「新宿／泊棟」の 2 つがあったとする。後者は「泊棟」が未知語推定されたものとする。
文字間接続確率は図 26（b）のように与えられているものとする。この例の場合、先に示した式（7）だけで計算すると、図 28（a）の候補（1）はスコアが 0.044 に、候補（2）はスコアが 0.040 になり、候補（2）が誤って選択されてしまう。

【0098】

これに対して先に示した式（11）の計算方法を採用すると、Th として 0.03 を与えた時、図 28（c）のように候補（1）は「宿泊」の「宿」と「泊」の間は単語に切れないのでこの位置に定数 Th を与えてスコアを計算することで、候補（1）のスコアとして 0.074 を得る。同様に候補（2）は「新」と「宿」の間、および「泊」と「棟」の間に定数 Th を与えて計算することで、候補（2）のスコアとして 0.100 を得る。これらを比較することで候補（1）

を選択することができる。

【0099】

式(7)を用いた分割が分割数が最小になるものを選択する傾向にあるのに対して、式(11)を用いた計算は、より細かく分割した解候補の選択をも可能とする。つまり、「衆議院議員」のような複合語が辞書にあったとしても、それを「衆議院」と「議員」という短い単位の単語に分解できる。単語分割装置の使用目的によっては、文書によって分割の粒度を制御する要求があるが、Thをパラメータとすることでこれを実現可能とする。

【0100】

式(11)は、スコアに和を取る場合の式(7)の計算に閾値に相当する値を導入するものであるが、式(8)の積の場合にも同様に閾値を導入することで、分割粒度を制御することが可能になる。

【0101】

また、これまでに述べたスコア計算方法は単語分割位置についての文字間接続確率と、単語分割されない位置についての定数値のみを用いたが、本発明では夫々の単語に何らかのスコアを与えることも提案する。単語へのスコアの与え方として、それが単語辞書記憶手段207に記憶されていたものであれば定数Uを、未知語推定手段210が推定した未知語であれば定数Vを与えるものとし、候補中の単語の集合をW、単語辞書記憶手段207に記憶されている単語の集合をDとすると、式(11)を拡張したスコアは次の式により記述できる。

【0102】

【数 14】

$$\arg \min_{S \subseteq N, W} \left(\sum_{i \in N} Q_i + \sum_{j \in W} R_j \right) \quad (14)$$

$$Q_i = \begin{cases} P_i, & i \in S \\ Th, & i \notin S \end{cases} \quad (15)$$

$$R_j = \begin{cases} U & j \in D \\ V & j \notin D \end{cases} \quad (16)$$

$$U < V \quad (17)$$

【0103】

この時、 $U < V$ とすることで未知語よりも辞書に載っている単語を優先して選択することを可能とする。よって、単語分割候補の中に含まれている未知語が少ないものが解として選択されるという効果を生む。

【0104】

未知語推定手段210の導入と、定数 Th の導入、および単語のスコアの導入により、未知語推定をしながら推定した未知語を選択するか否か判断しながら単語分割をすることになる。未知語推定手段は図23のステップ2302で辞書に無い長さ $n \sim m$ 文字の任意の文字列を未知語候補とした。この未知語を文字間接続確率によって選択することは、未知語の部分は本発明の第1または第2の実施の形態のように確率値だけで分割することと等価になる。よって、辞書による単語分割と確率による単語分割の統合が可能となる。

【0105】

従来の技術では未知語の推定には、「漢字と平仮名の境界は単語の分割点になりやすい」といった経験による知識を用いていたが、本発明においては図23のステップ2302において条件範囲内の全ての文字列を未知語として扱う。しかし文字間接続確率の計算をすることでその中から正解を選択することができる。これは例えば漢字と平仮名にまたがった未知語も推定できることを意味している

【0106】

以上のように、本発明の第4の実施の形態においては、分割処理対象文書自身から文字間接続確率を計算しておき、単語辞書と未知語推定手段を使って分割解の候補を作成するが、複数の解候補から一つを選択するには文字間接続確率が最小になるものを選択する。辞書にない単語については未知語を推定するが、未知語を選択するかどうかは確率値（あるいは計算されたスコア）に従うことで、未知語周辺は確率値により分割を行う。よって分割候補選択のための知識の学習に、あらかじめ人手で作成された分割正解文書を大量に用意する必要がないことからコストを抑えることが可能であり、分割対象の文書から確率の形で知識を自動学習することが可能で、文書分野に適合した学習ができる点で合理的であり、その実用的効果は大きい。さらに辞書にない未知語は確率値から分割することも可能であり、その実用的効果は大きい。

【0107】

【発明の効果】

以上のように、本発明では、処理対象文書の中から文字間接続確率を計算し、その確率を処理対象文書にあてはめることで単語分割できる場所を発見して分割するものであり、これにより、辞書を使わずにテキストを単語に分割するという効果を奏するものである。よって本発明は単語の分割のために辞書を用意する必要がないことから、日々生まれ続ける新しい語や語法のために、辞書を整備したり各種パラメータを整備する必要もない。辞書を持たないことから、本方式により実現されたプログラムを格納した記録媒体は非常に小さくて済む。同時にパーソナルコンピュータなどの処理能力に限界のある環境下においても機能する。

【0108】

また、辞書を用いた単語分割においても文字間接続確率を用いることで、複数の分割解候補の中から一つを選択することが可能となる。

【0109】

さらに、辞書を用いた単語分割において辞書に載っていない単語が出現しても、未知語推定と文字間接続確率を併用することで文字列を単語単位に分割するこ

とが可能となる。

【図面の簡単な説明】

【図 1】

本発明の第 1 の実施の形態における単語分割方式の動作を示すフローチャート

【図 2】

本発明の第 1 の実施の形態における単語分割装置の構成を示すブロック図

【図 3】

本発明の第 1 の実施の形態における文字間接続確率の計算手順を示すフローチャート

【図 4】

本発明の第 1 の実施の形態における単語分割計算例を示す概念図

【図 5】

本発明の第 1 の実施の形態における分割過程の計算手順を示すフローチャート

【図 6】

本発明の第 1 の実施の形態における単語分割計算例を示す概念図

【図 7】

本発明の第 1 の実施の形態における単語分割方式で新聞記事データの文字間接続確率を計算した例の一部を示す数値図

【図 8】

本発明の第 1 の実施の形態における単語分割方式で新聞記事データの分割をした例を示す概念図

【図 9】

本発明の第 2 の実施の形態における文字間接続確率の計算手順について、 n と m が違う場合の計算手順を示すフローチャート

【図 10】

本発明の第 2 の実施の形態における文字間接続確率の計算手順について、 n と m が同じ場合の計算手順を示すフローチャート

【図 11】

本発明の第 2 の実施の形態における分割過程の計算手順を示すフローチャート

【図 1 2】

本発明の第 2 の実施の形態における単語分割計算例を示す概念図

【図 1 3】

本発明における閾値と平均単語長の関係を示す概念図

【図 1 4】

本発明の第 2 の実施の形態における式 (5) の関係を示す模式図

【図 1 5】

本発明の第 3 の実施の形態における単語分割装置の構成を示すブロック図

【図 1 6】

本発明の第 3 の実施の形態における単語分割方式の動作を示すフローチャート

【図 1 7】

本発明の第 3 の実施の形態における単語分割過程の例を示す概念図

【図 1 8】

本発明の第 3 の実施の形態における単語分割過程の例を示す概念図

【図 1 9】

本発明の第 3 の実施の形態における単語分割の詳細動作を示すフローチャート

【図 2 0】

本発明の第 3 の実施の形態における辞書利用過程を示す概念図

【図 2 1】

本発明の第 4 の実施の形態における単語分割装置の構成を示すブロック図

【図 2 2】

本発明の第 4 の実施の形態における単語分割方式の動作を示すフローチャート

【図 2 3】

本発明の第 4 の実施の形態における単語分割の詳細動作を示すフローチャート

【図 2 4】

本発明の第 4 の実施の形態における単語分割の辞書内容の例を示す概念図

【図 2 5】

本発明の第 4 の実施の形態における辞書利用と未知語推定の過程を示す概念図

【図 2 6】

本発明の第 4 の実施の形態における単語分割過程の例を示す概念図

【図 2 7】

本発明の第 4 の実施の形態における分割解の候補の選択過程の例を示す概念図

【図 2 8】

本発明の第 4 の実施の形態における分割解の候補の選択方法の例を示す概念図

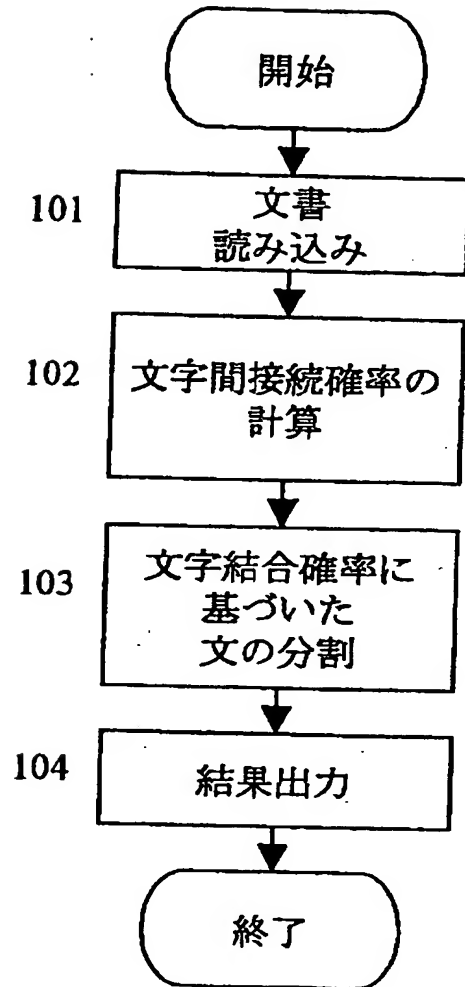
【符号の説明】

- 2 0 1 文書入力手段
- 2 0 2 文書データ蓄積手段
- 2 0 3 文字間接続確率計算手段
- 2 0 4 確率テーブル記憶手段
- 2 0 5 文字列分割手段
- 2 0 6 文書出力手段
- 2 0 7 単語辞書記憶手段
- 2 0 8 分割解候補作成手段
- 2 0 9 解選択手段
- 2 1 0 未知語推定手段

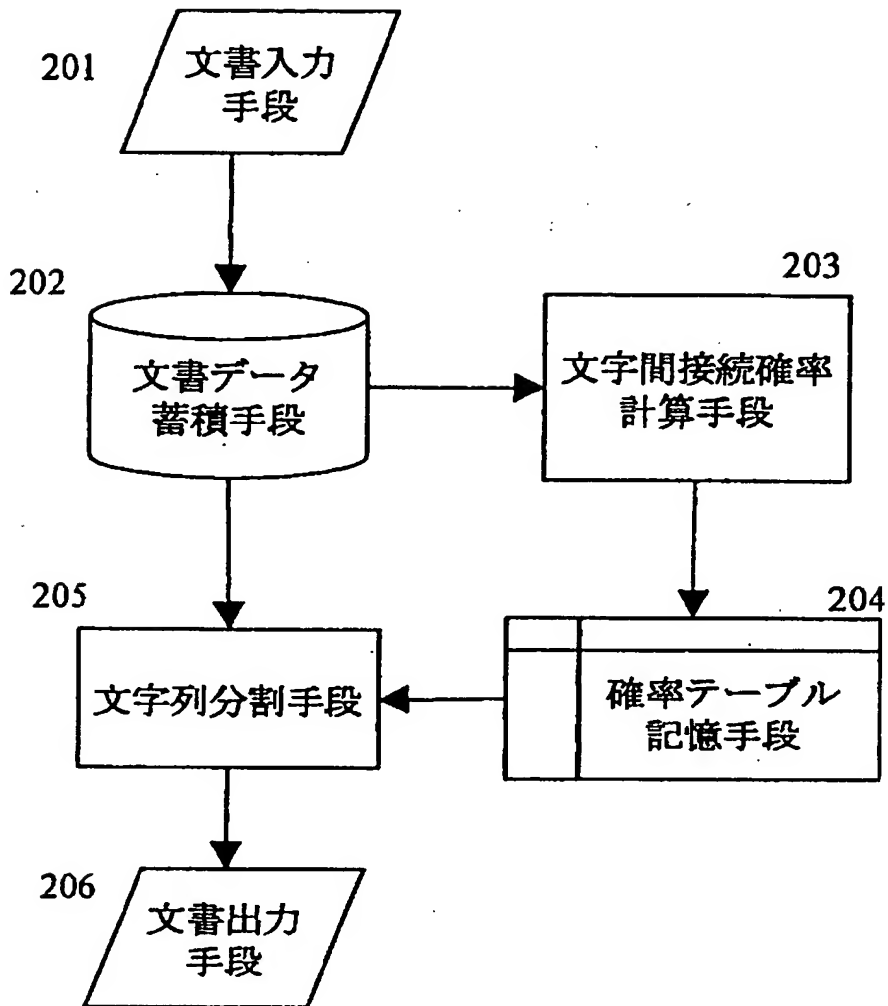
【書類名】

図面

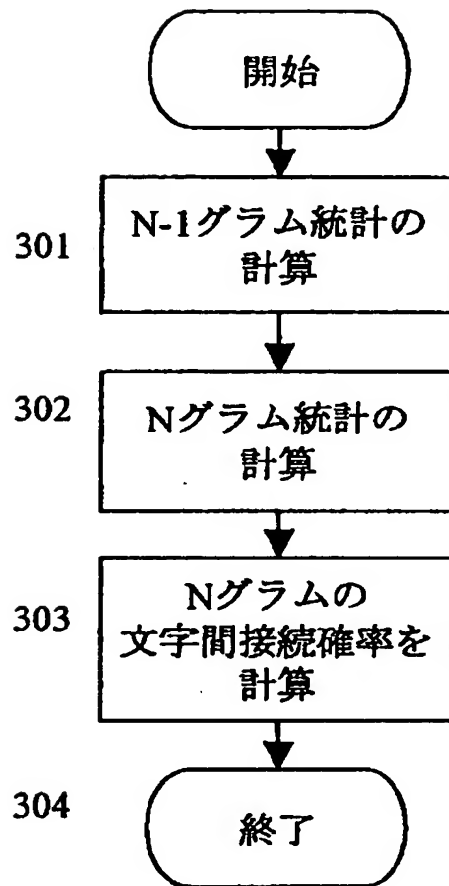
【図 1】



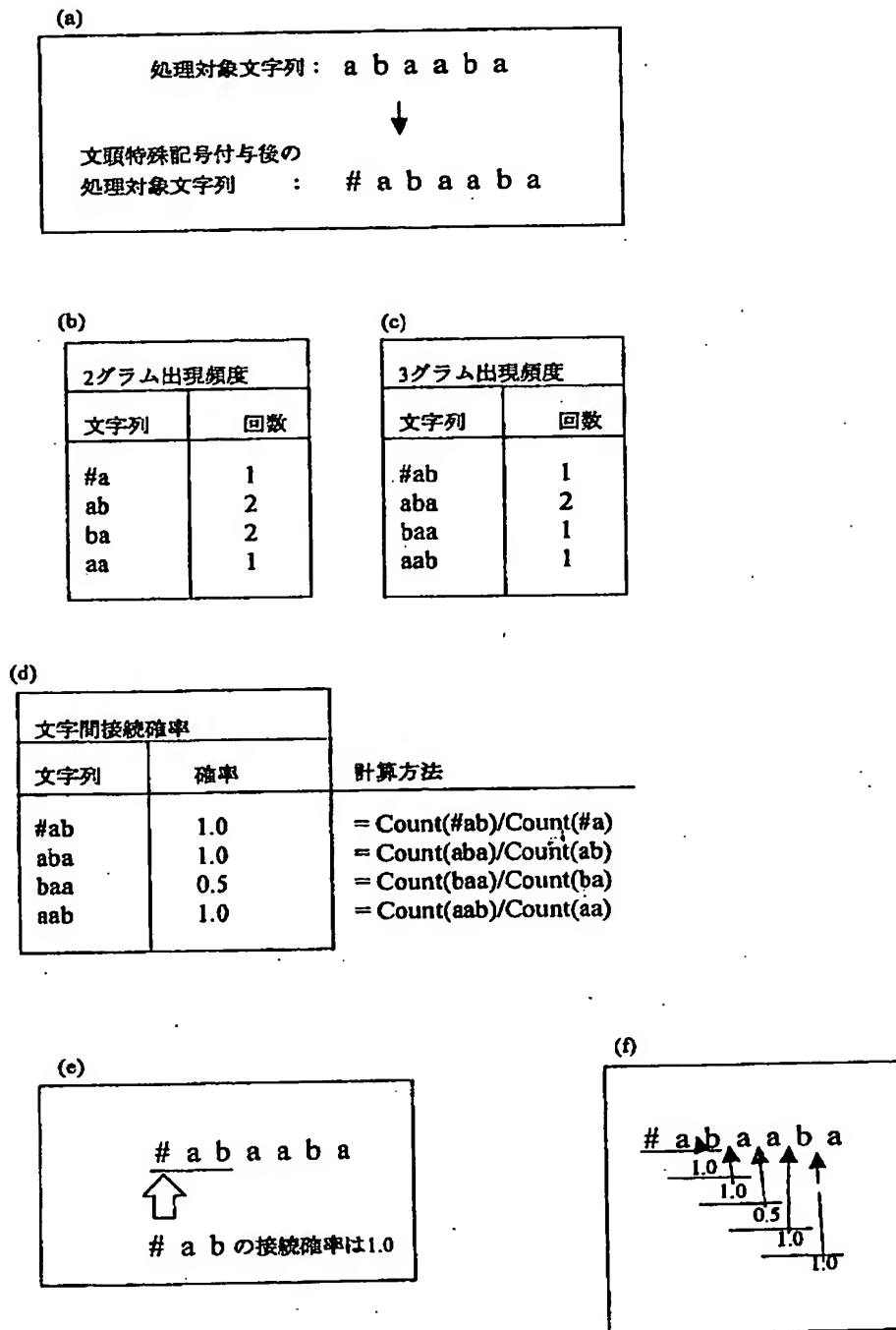
【図 2】



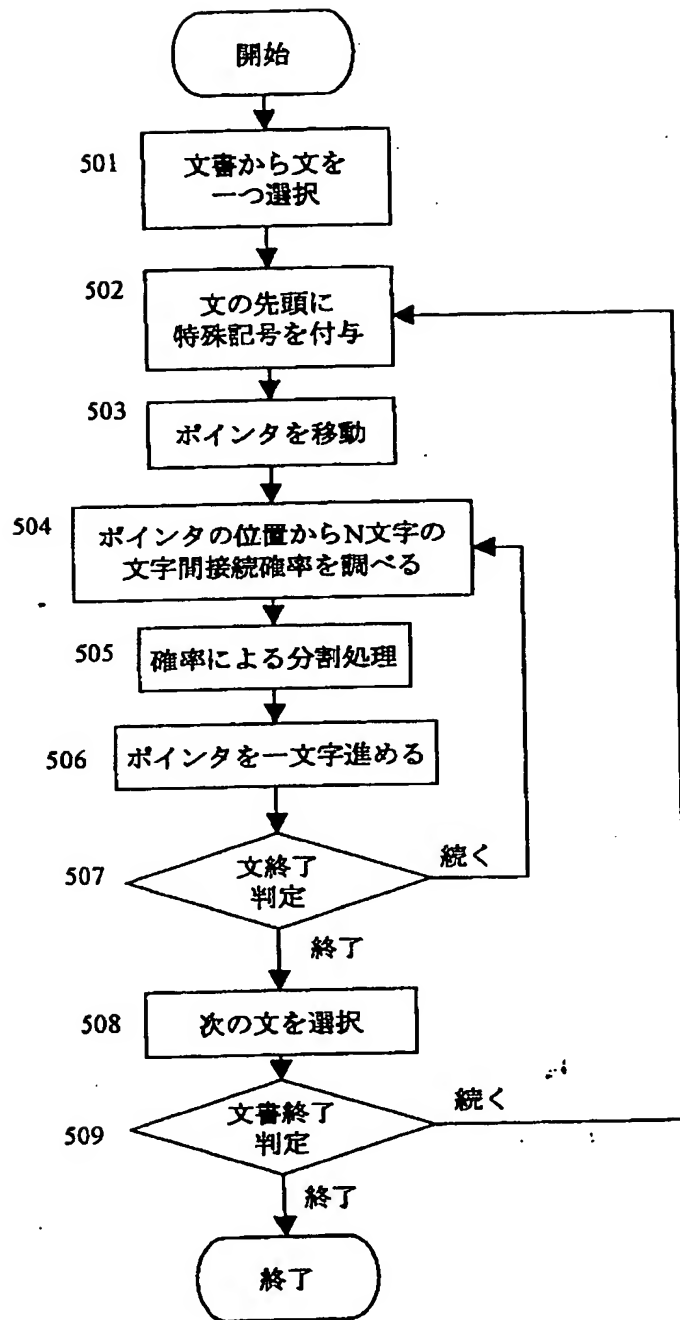
【図 3】



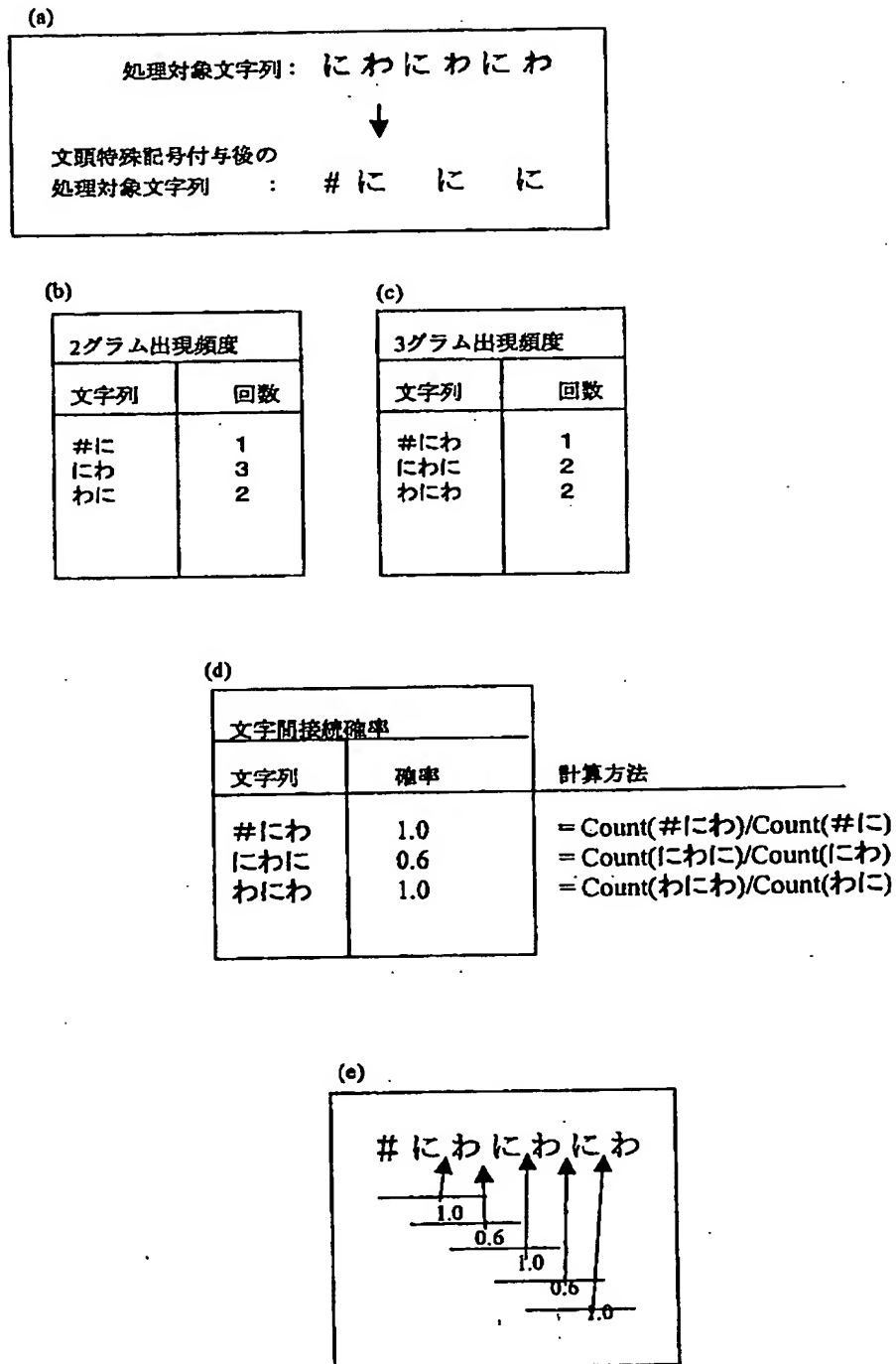
【図 4】



【図 5】



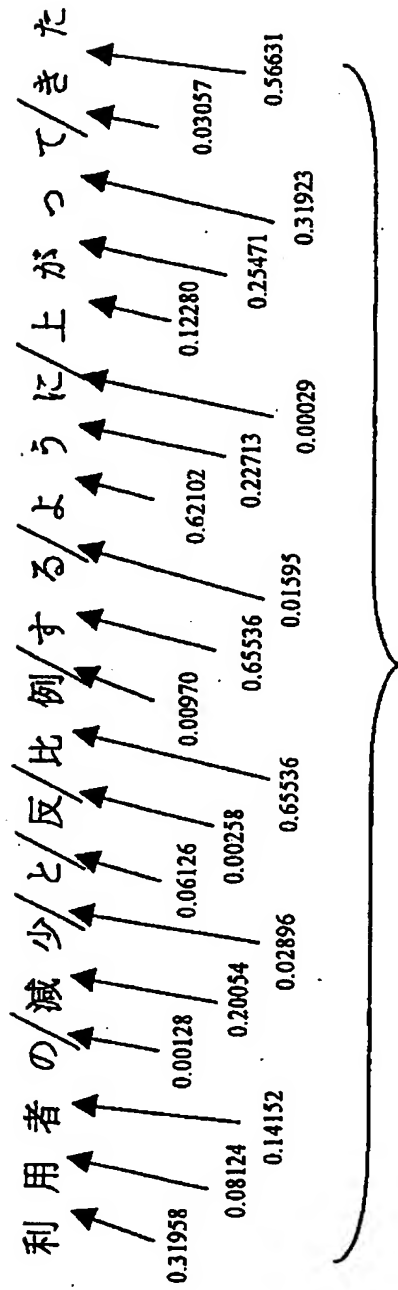
【図 6】



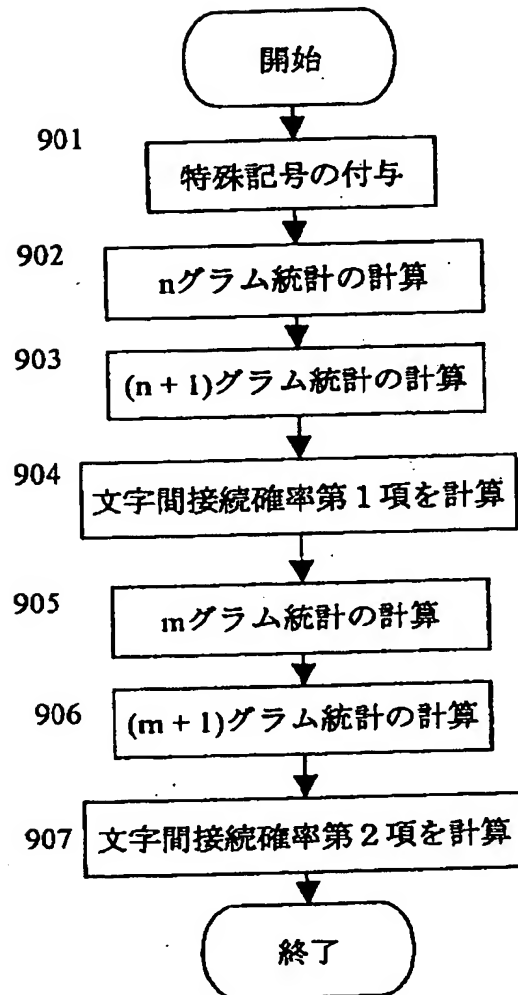
【図 7】

文字間接続確率	
文字列	確率
ようち	0.00016
ようっ	0.00011
ようつ	0.00005
ようて	0.00011
ようで	0.01938
ようと	0.08398
ような	0.13665
ように	0.22713
よくあ	0.01593
よくい	0.00612
よくえ	0.00024
よくお	0.00269
よくか	0.00171

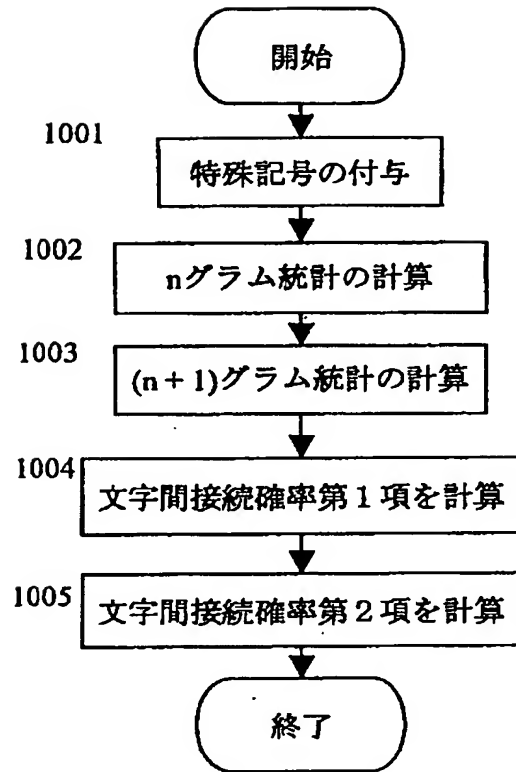
【図8】



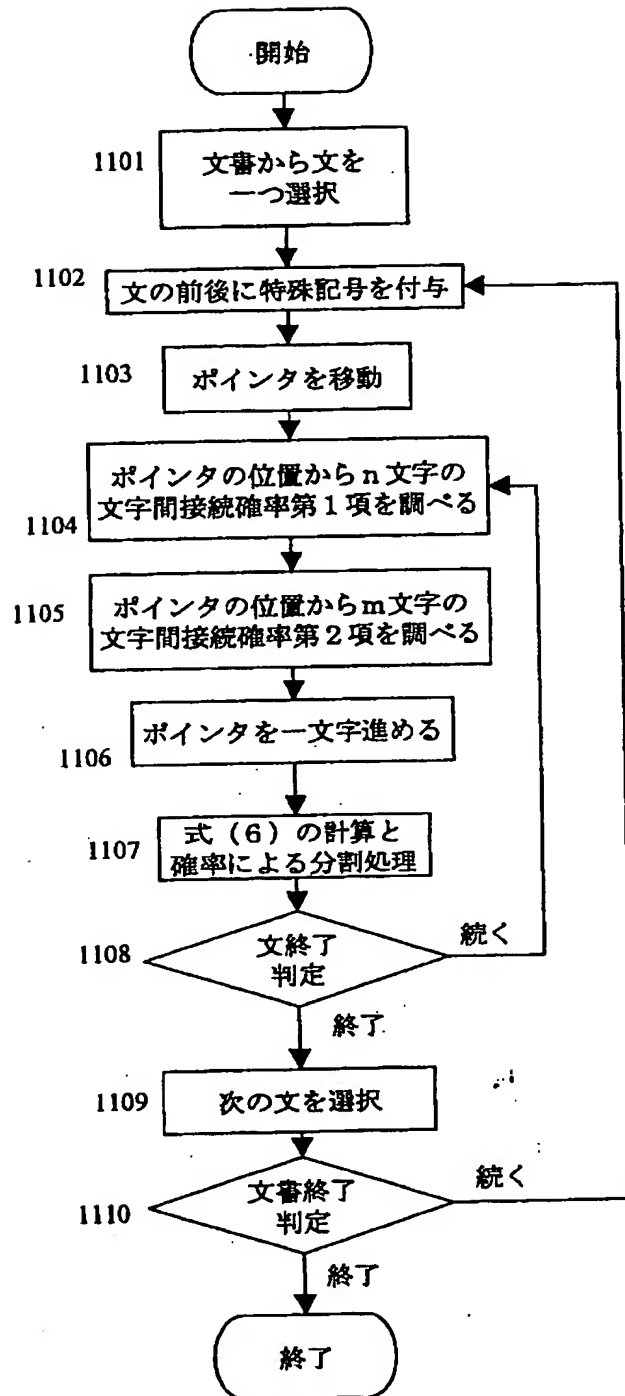
【図 9】



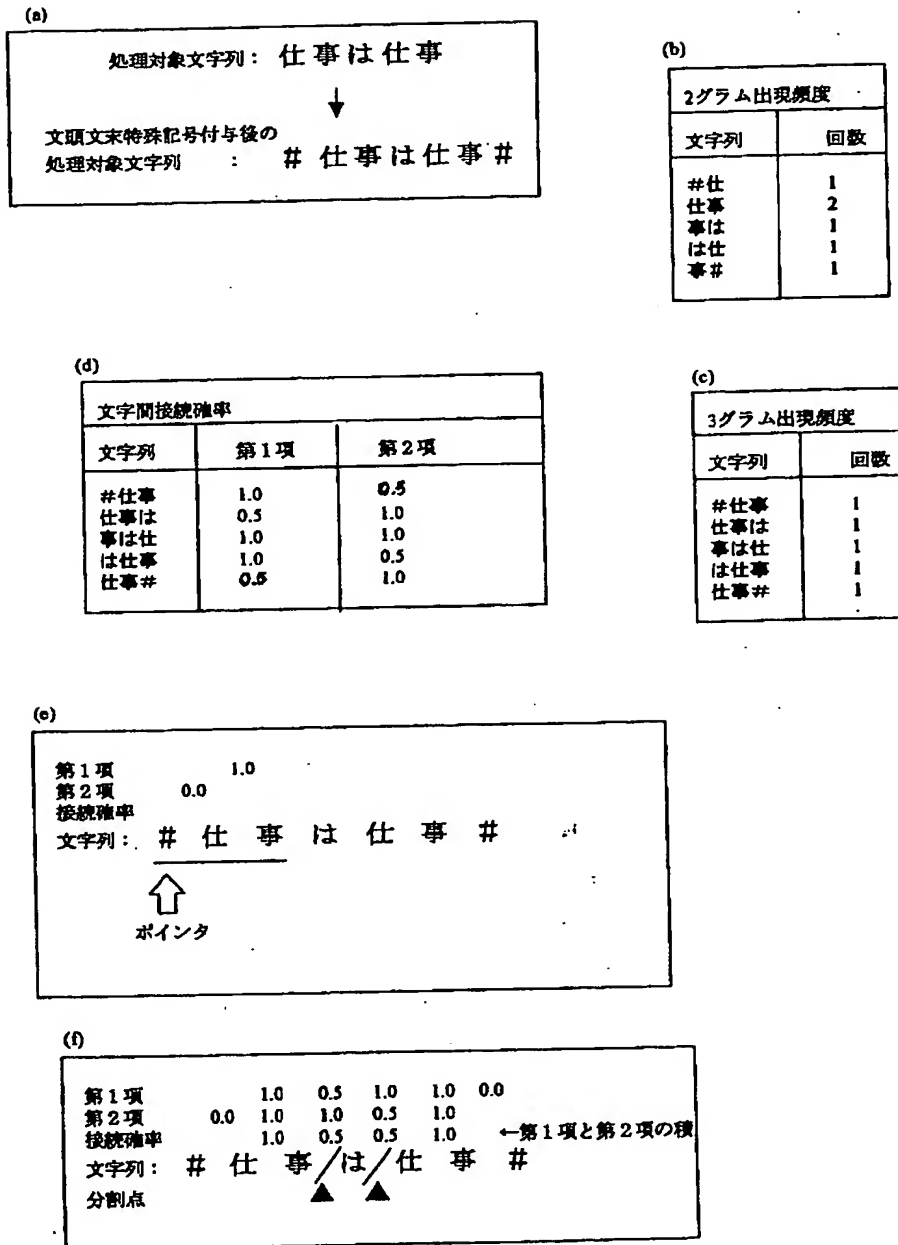
【図 1 0】



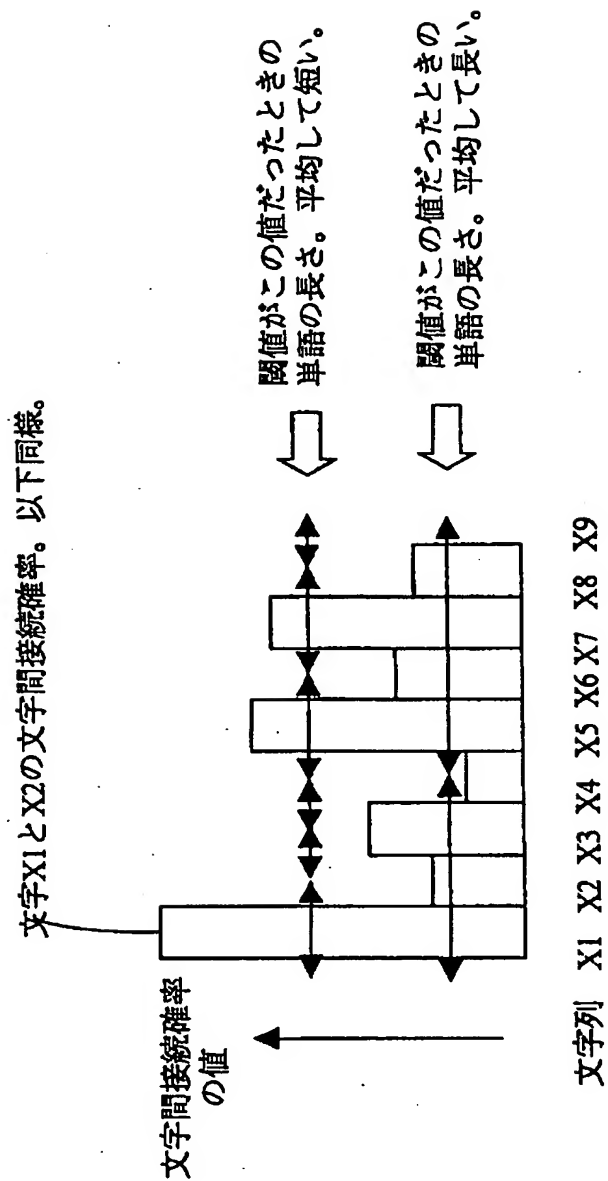
【図 11】



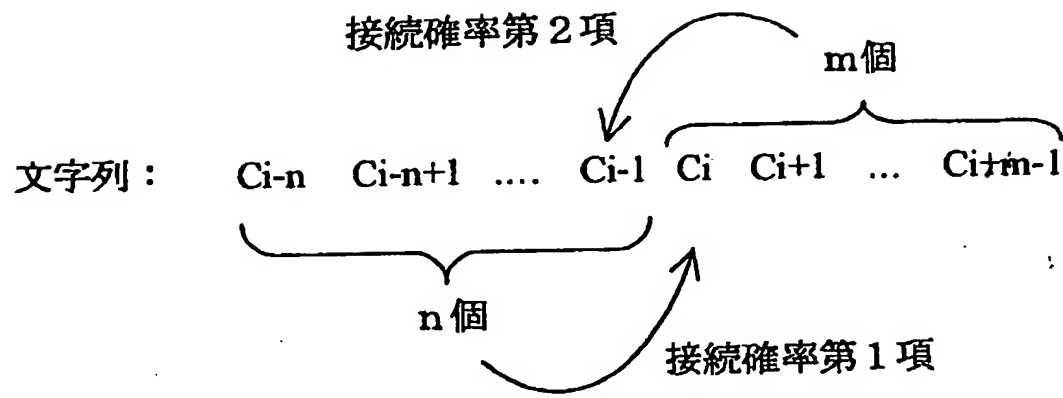
【図 12】



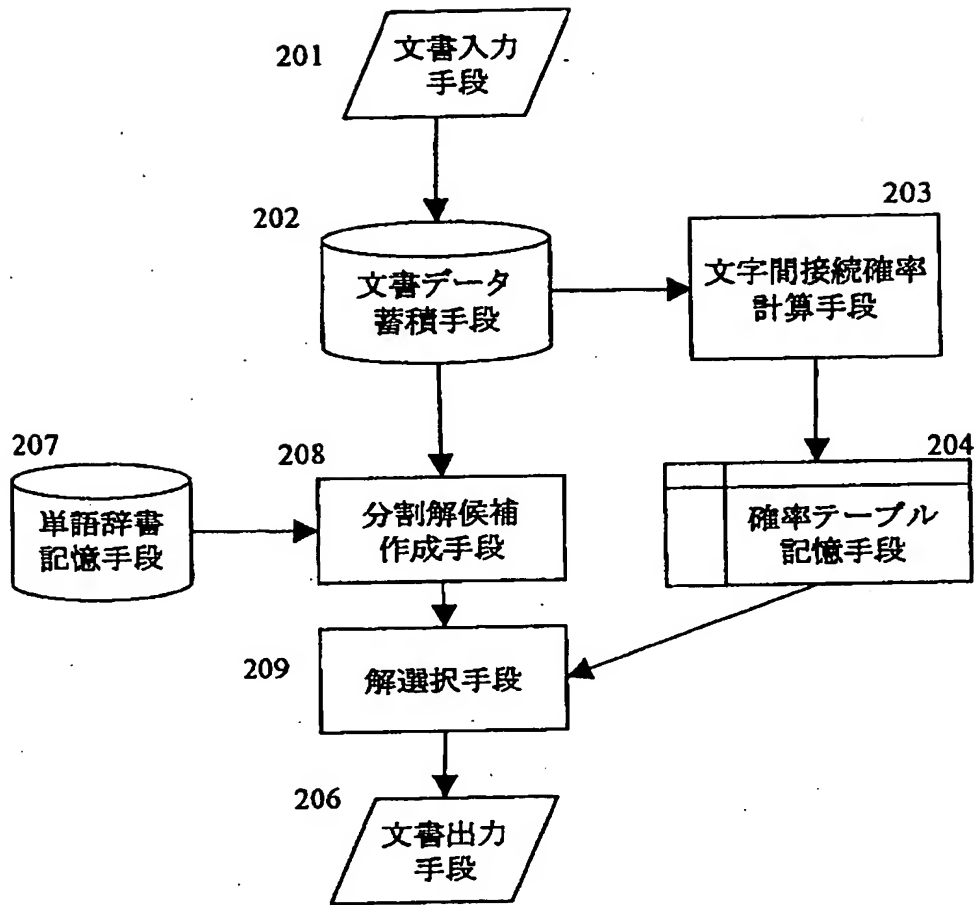
【図 1 3】



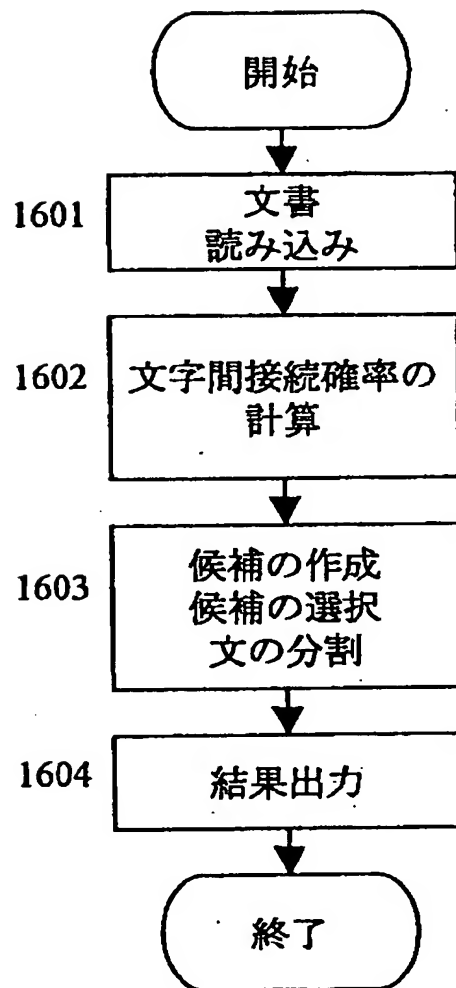
【図 14】



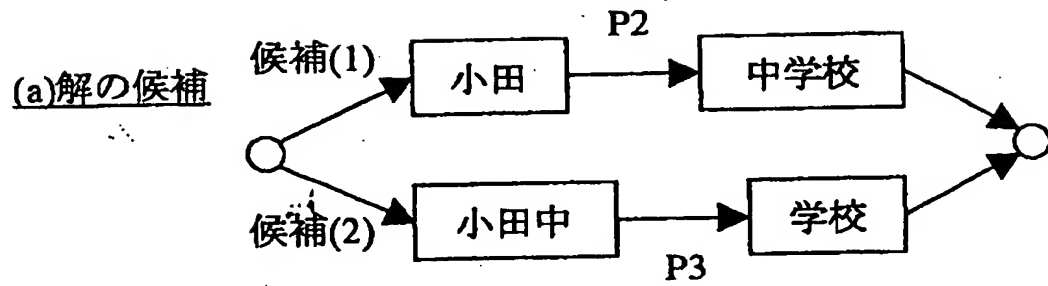
【図 15】



【図 16】



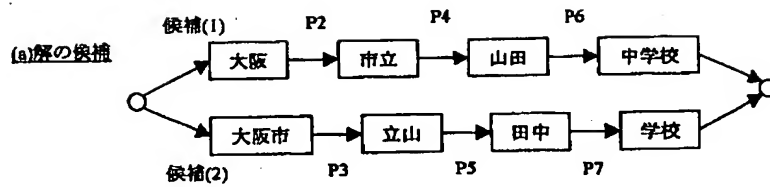
【図17】



(b) 文字間の接続確率表

小	>	P1 = 0.019
田	>	P2 = 0.007
中	>	P3 = 0.088
学	>	P4 = 0.439
校		

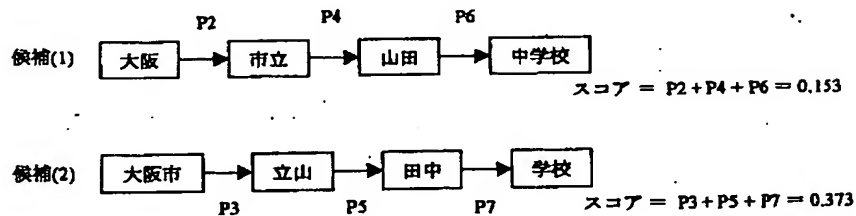
【図 18】



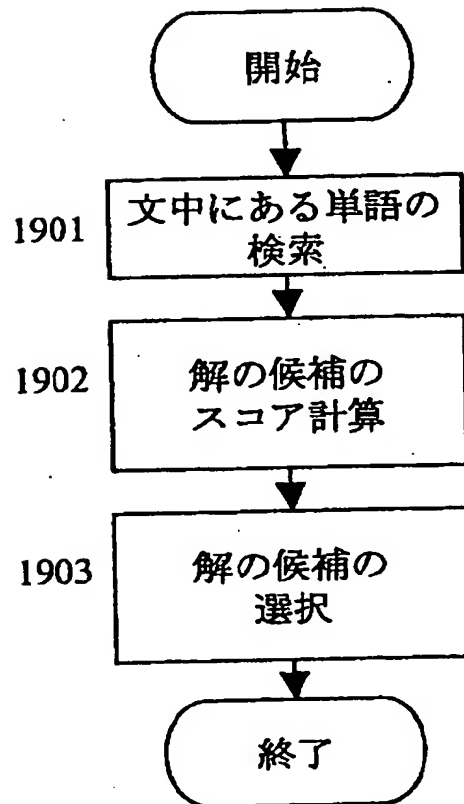
(b) 文字間の接続確率表

大	>	P1 = 0.599
阪	>	P2 = 0.141
市	>	P3 = 0.174
立	>	P4 = 0.006
山	>	P5 = 0.111
田	>	P6 = 0.006
中	>	P7 = 0.088
学	>	P8 = 0.439
校		

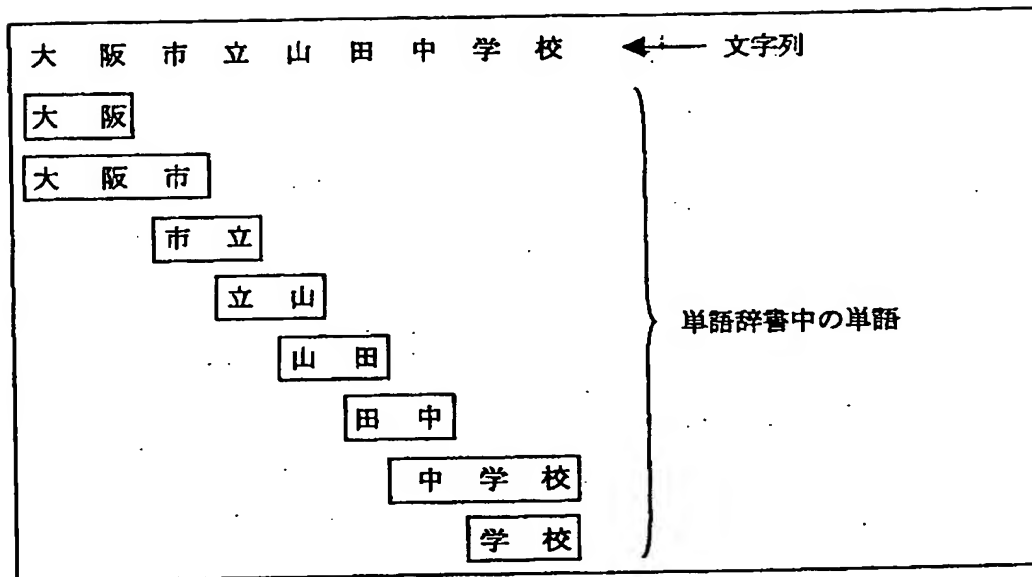
(c) 解の候補の比較



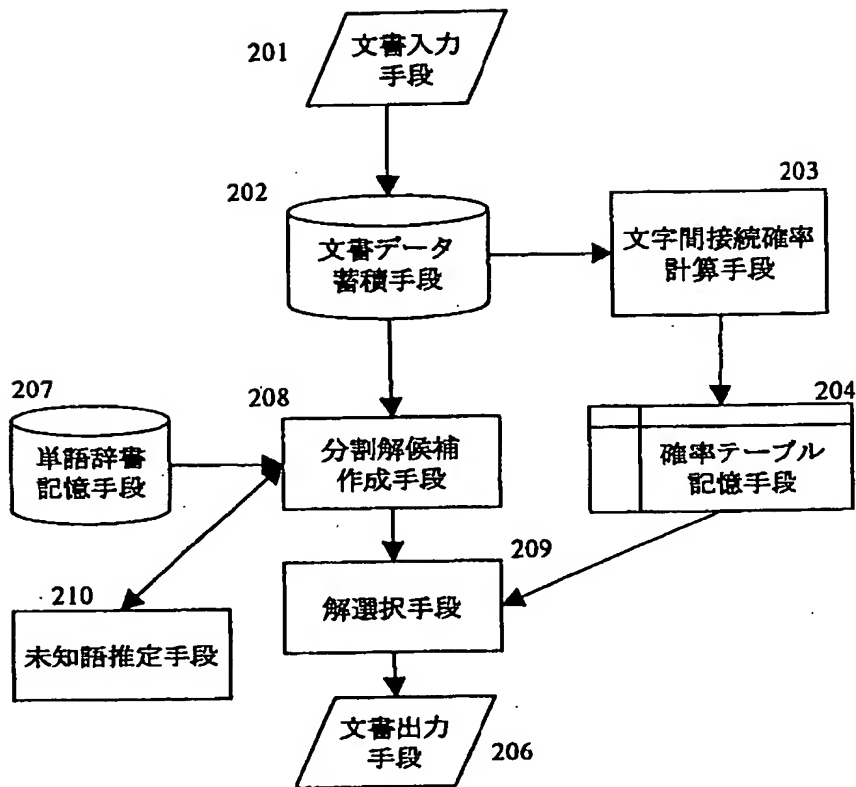
【図 19】



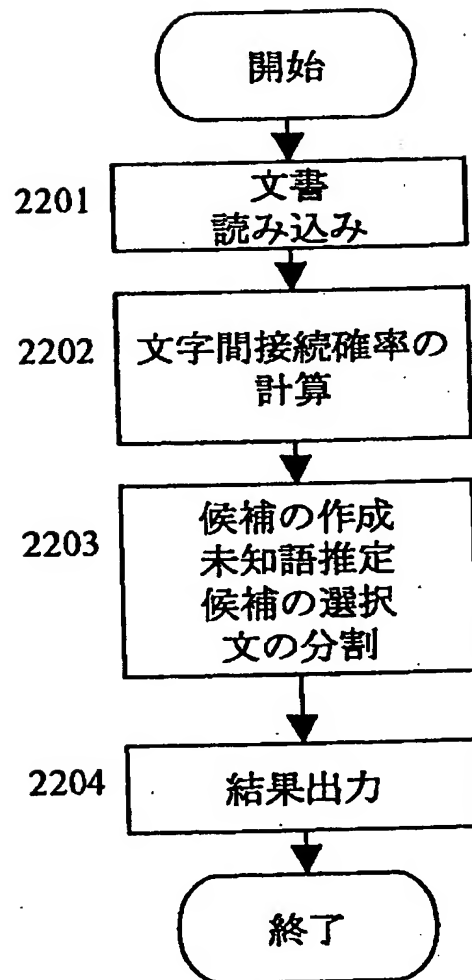
【図20】



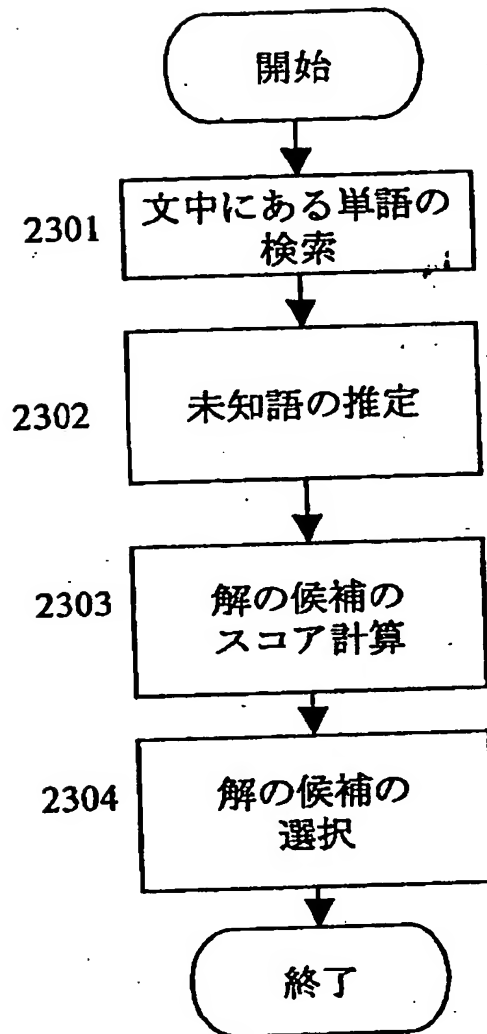
【図 21】



【図 2 2】

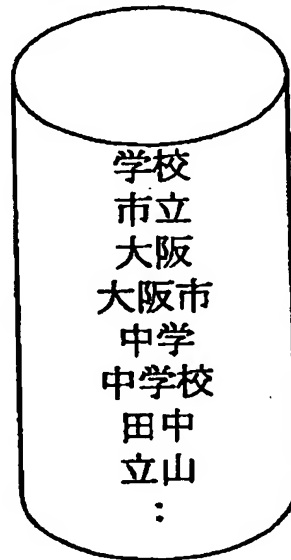


【図 23】

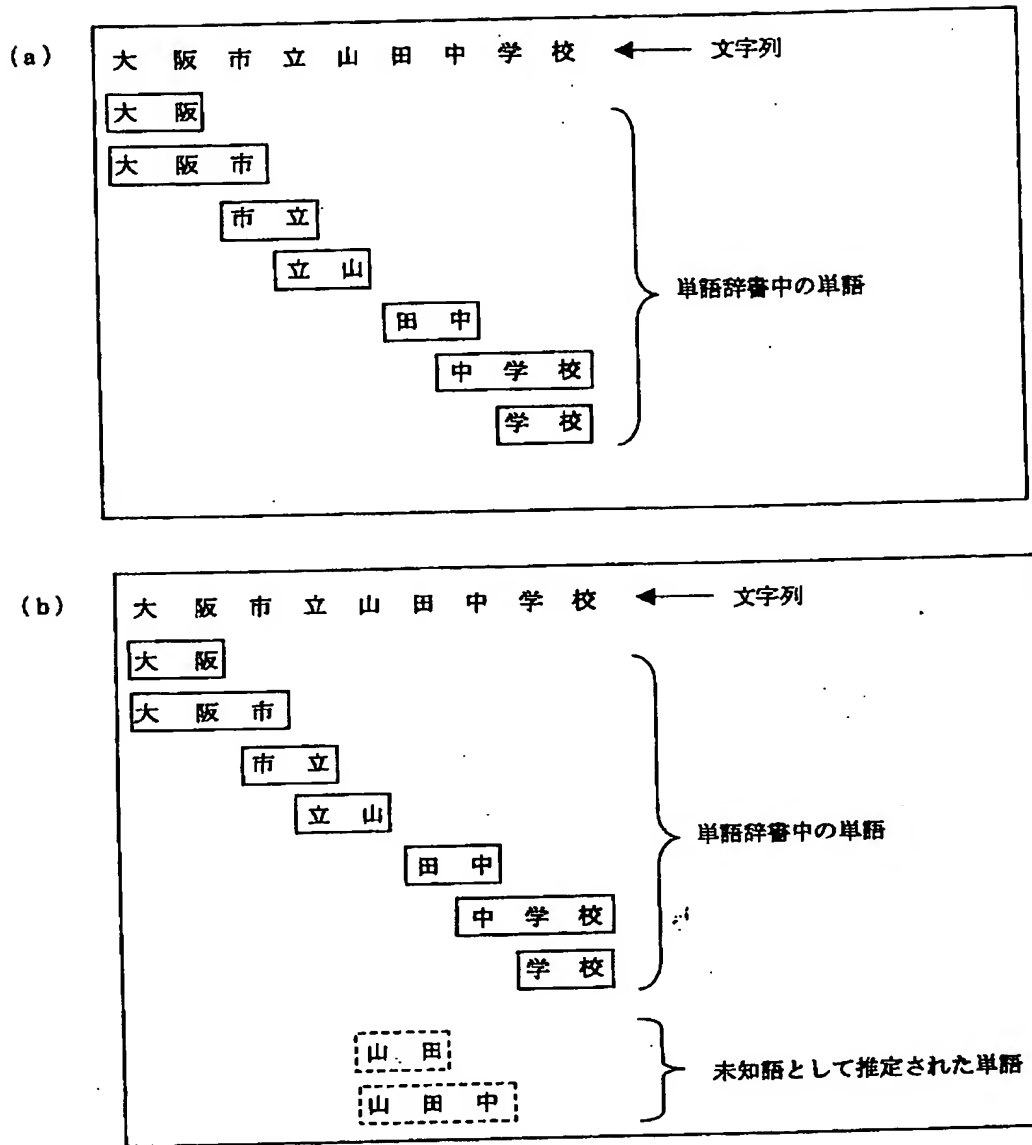


【図24】

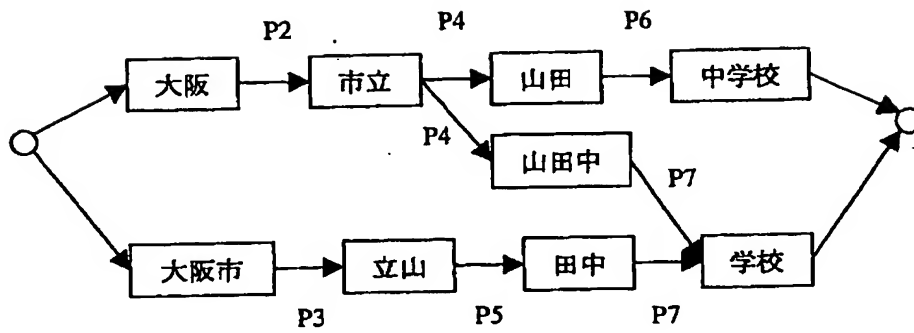
単語辞書記憶手段207



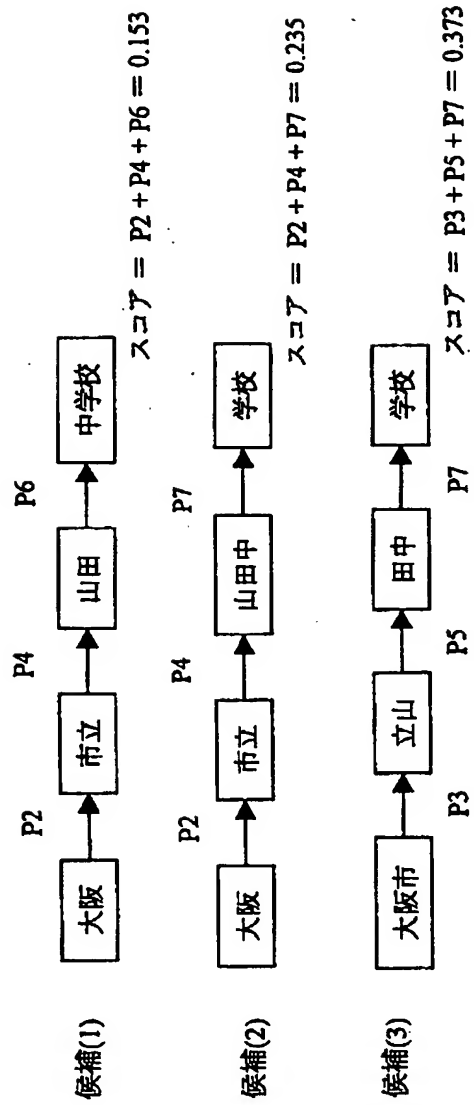
【図25】



【図 26】

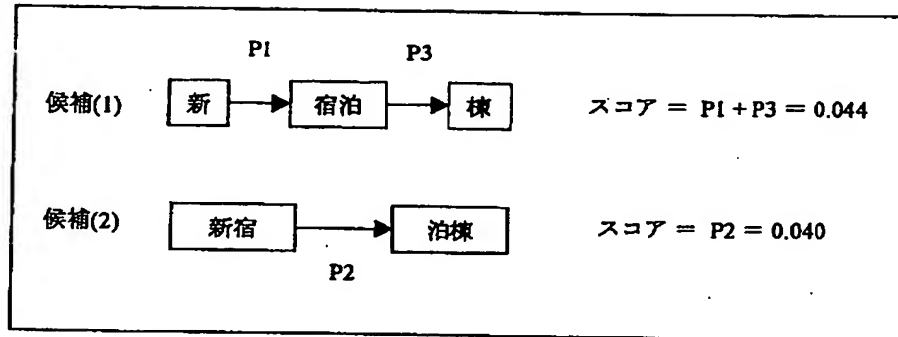


【図27】



【図 28】

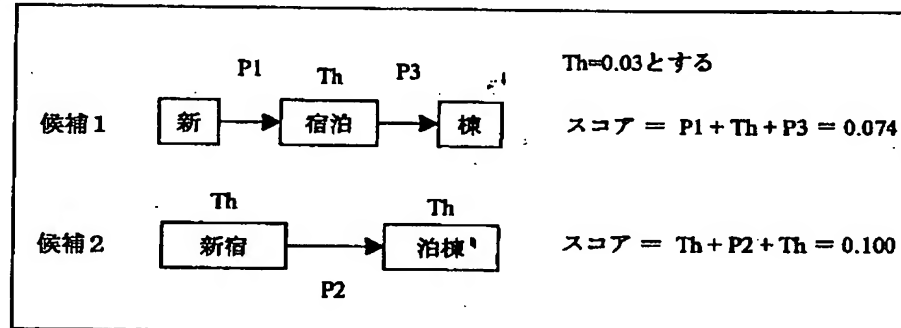
(a) 解の候補の比較



(b) 文字間の接続確率表

新	> $P1 = 0.016$
宿	> $P2 = 0.040$
泊	> $P3 = 0.028$
棟	

(c) 解の候補の比較



【書類名】 要約書

【要約】

【課題】 電子計算機を利用した自然言語処理システムにおいて、辞書や学習用単語分割済文を必要としない単語分割方式を実現することを目的とする。

【解決手段】 入力された単語分割されていない文書から、文字結合度としての文字間接続確率を統計的に計算し、テーブルに記録する。この文字間接続確率を用いて入力された文書を調べ、文字間接続確率の低い部分で文書を分割し、出力する。

【選択図】 図 1

出 願 人 履 歴 情 報

識別番号 [000005821]

1. 変更年月日 1990年 8月28日

[変更理由] 新規登録

住 所 大阪府門真市大字門真1006番地
氏 名 松下電器産業株式会社